

Fairness Analysis Report

Information about test dataset

Grouping by sex

category	size	number patients with event
F	1565	304
M	2834	734

Grouping by age_group

category	size	number patients with event
<50	467	59
50-65	1292	305
65-75	1335	379
75-85	1056	258
>85	249	37

Grouping by APACHE_group

category	size	number patients with event
Cardiovascular	1307	464
Neurological	1490	203
Gastrointestinal	545	98
Respiratory	332	128
Other	283	53
Trauma	300	67
Metabolic	111	23

Grouping by surgical_status

category	size	number patients with event
Surgical	2186	464
Non-surgical	2128	548

Grouping by admission_type

category	size	number patients with event
emergency	2743	719
elective	1591	299

Model Performance Analysis

Goal: Comparing the model performance across cohorts of patients

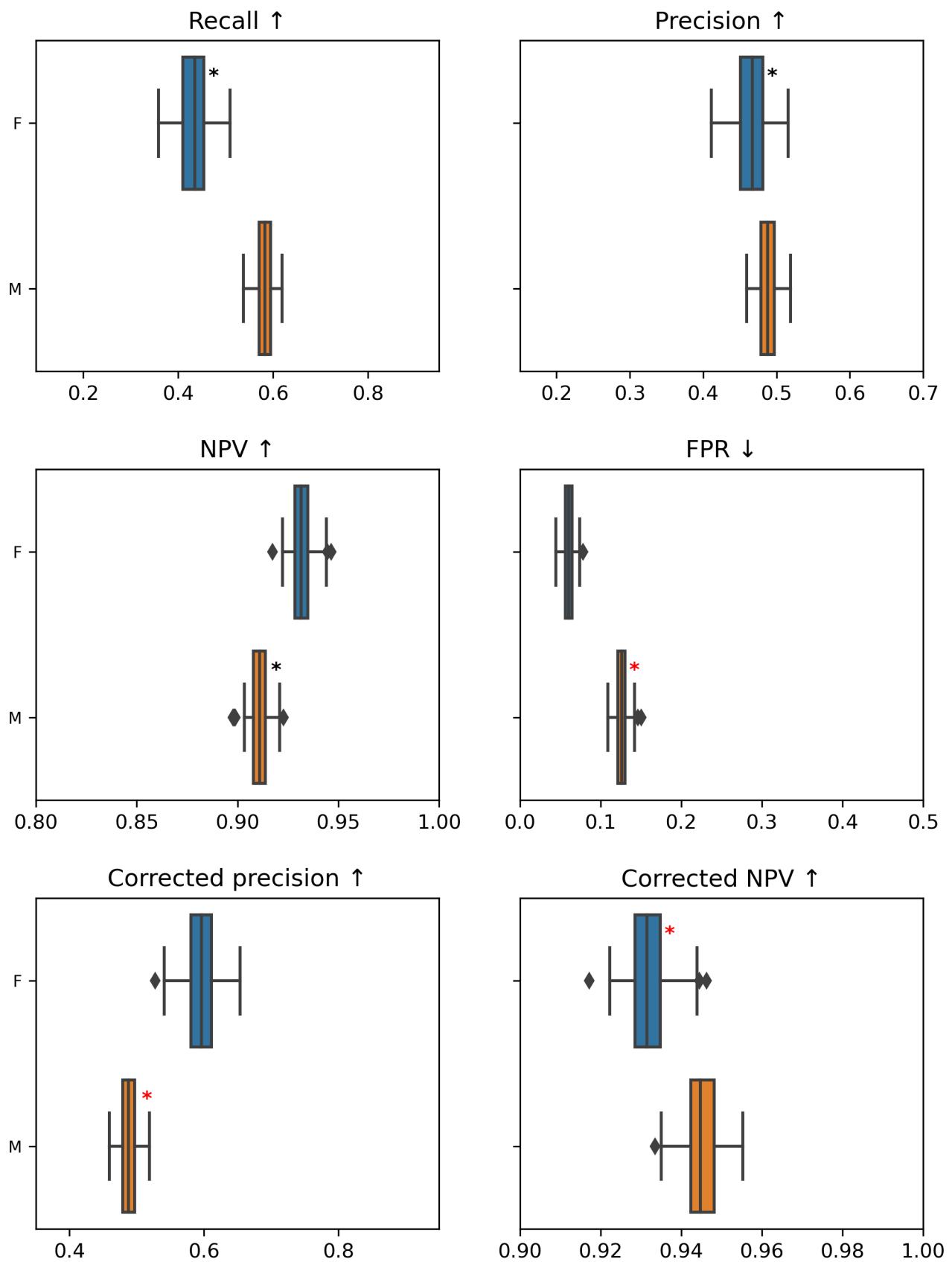
Binary metrics computed with a threshold on score of 0.351.

In the tables presenting the results of the statistical analysis, we present only metrics and groups with a significant p-value (smaller than 0.001/nb_comparison) and whose delta is bigger than 0.

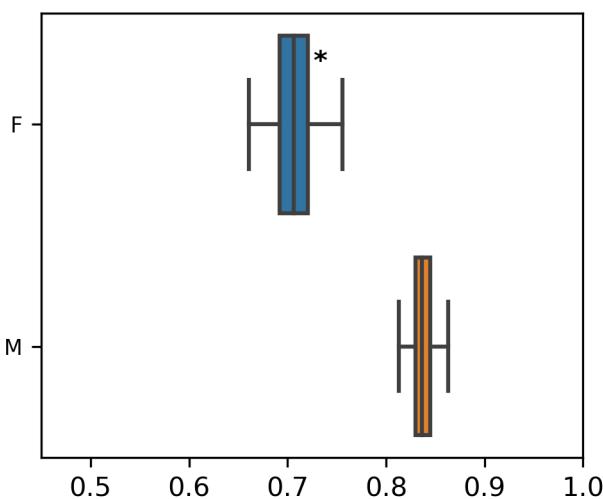
Top 3 categories with biggest deltas in terms of performance metric discrepancy

Metric	Cat 1	Cat 2	Cat 3
Recall ↑	Neurological (0.311)	Metabolic (0.209)	<50 (0.191)
Precision ↑	>85 (0.188)	<50 (0.138)	Neurological (0.085)
NPV ↑	Cardiovascular (0.087)	Respiratory (0.043)	65-75 (0.028)
FPR ↓	Respiratory (0.312)	Cardiovascular (0.166)	M (0.065)
Corrected precision ↑	M (0.108)	Respiratory (0.06)	Cardiovascular (0.058)
Corrected NPV ↑	F (0.013)	Surgical (0.013)	elective (0.007)
Event-based recall ↑	Neurological (0.296)	<50 (0.213)	Metabolic (0.15)
Avg. score on positive class	Neurological (0.139)	Metabolic (0.087)	<50 (0.081)
Avg. score on negative class	Respiratory (0.178)	Cardiovascular (0.115)	M (0.048)
AUROC ↑	Cardiovascular (0.074)	Respiratory (0.057)	75-85 (0.034)
AUPRC ↑	>85 (0.236)	Neurological (0.203)	<50 (0.166)
Corrected AUPRC ↑	>85 (0.046)	emergency (0.044)	M (0.031)
Event-based AUPRC ↑	>85 (0.351)	Metabolic (0.207)	Neurological (0.179)
Corrected event-based AUPRC ↑	>85 (0.157)	emergency (0.077)	Respiratory (0.058)

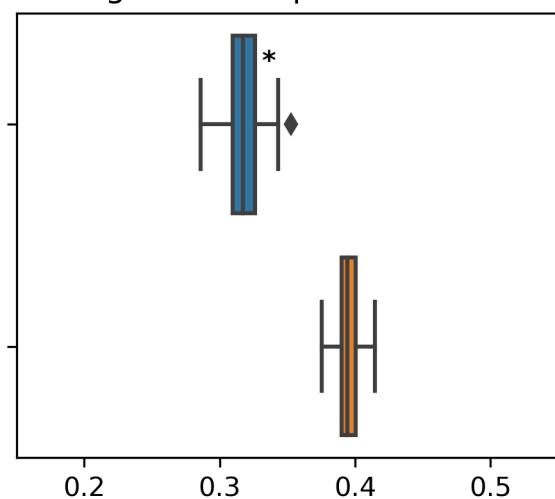
Grouping by sex



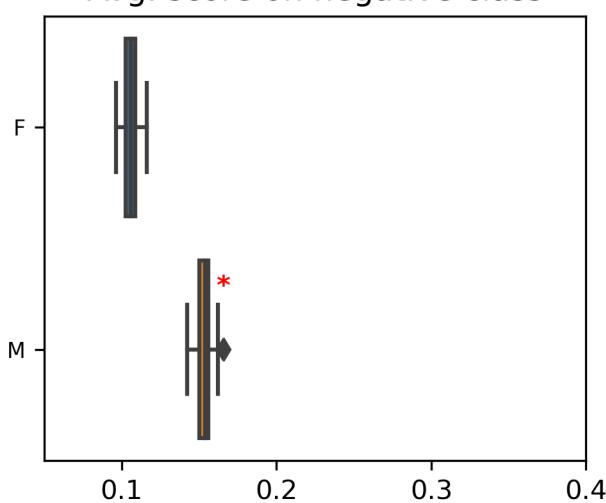
Event-based recall ↑



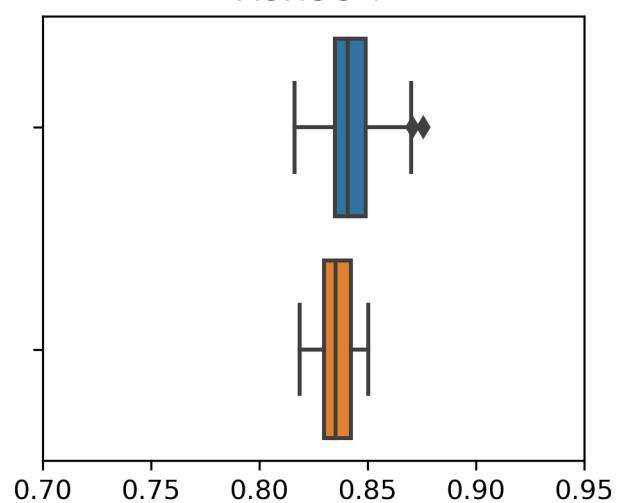
Avg. score on positive class



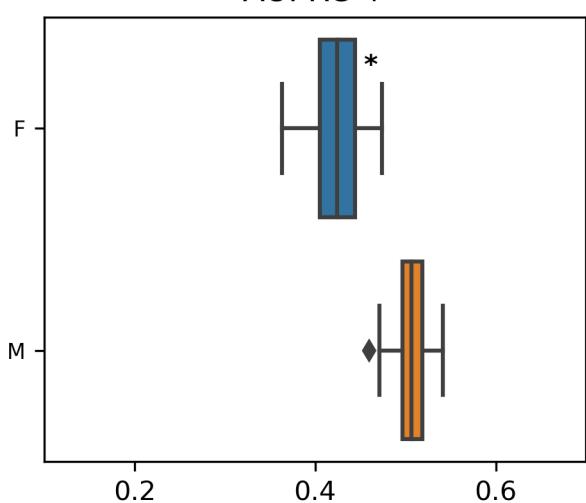
Avg. score on negative class



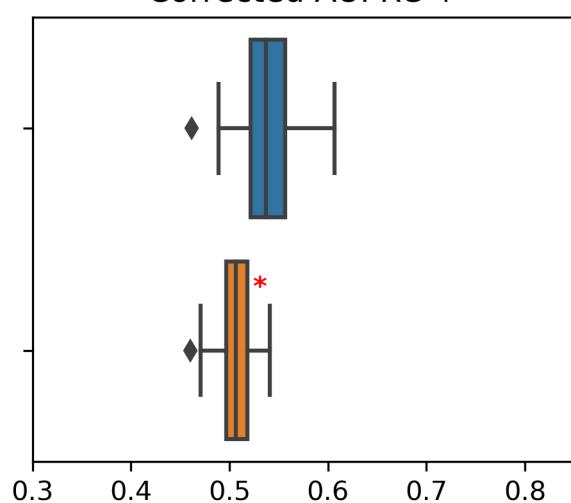
AUROC ↑

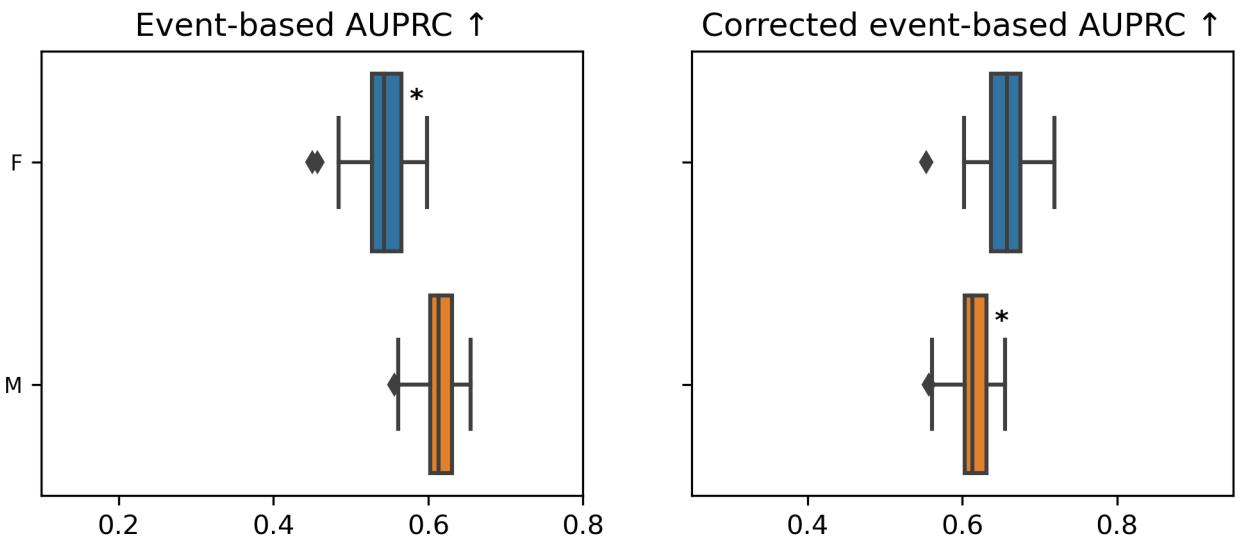


AUPRC ↑

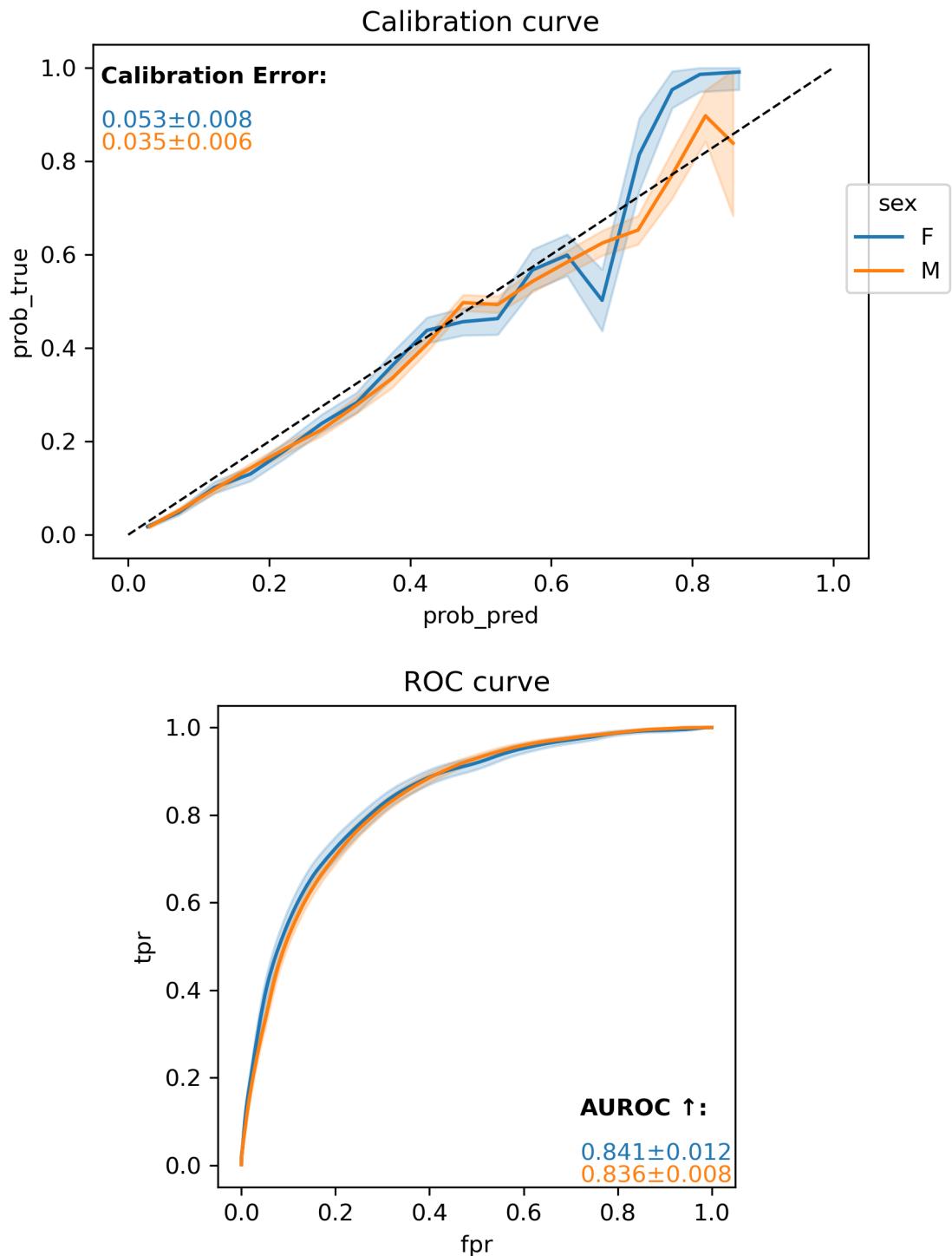


Corrected AUPRC ↑

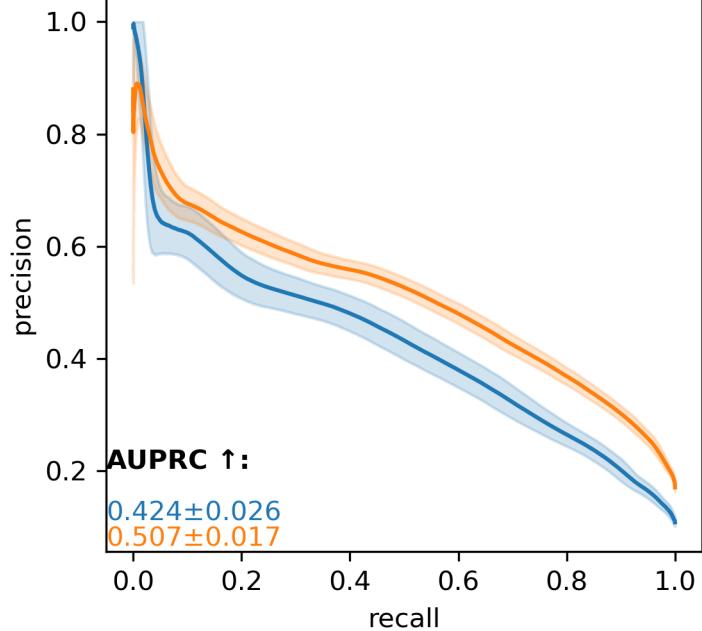




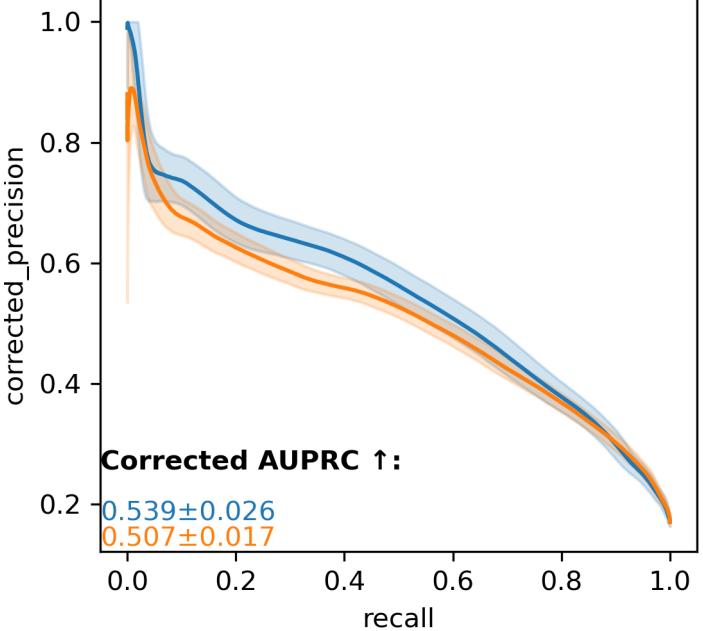
Metric	Group with worst dist.	P-value	Delta
Recall ↑	F	1.28e-34	0.148
Precision ↑	F	1.31e-12	0.021
NPV ↑	M	1.73e-34	0.021
FPR ↓	M	1.28e-34	0.065
Corrected precision ↑	M	1.28e-34	0.108
Corrected NPV ↑	F	1.93e-30	0.013
Event-based recall ↑	F	1.28e-34	0.13
Avg. score on positive class	F	1.28e-34	0.077
Avg. score on negative class	M	1.28e-34	0.048
AUPRC ↑	F	1.78e-34	0.082
Corrected AUPRC ↑	M	1.70e-18	0.031
Event-based AUPRC ↑	F	9.24e-32	0.07
Corrected event-based AUPRC ↑	M	1.48e-20	0.044



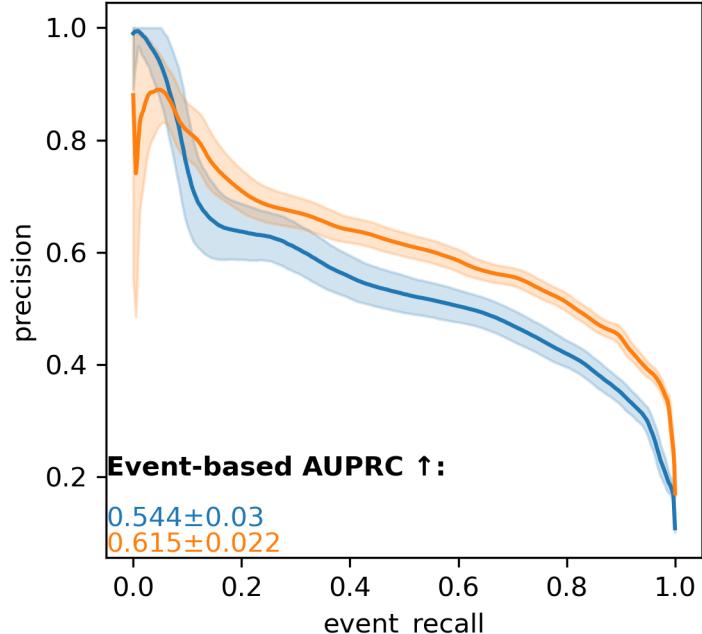
Precision / recall curve



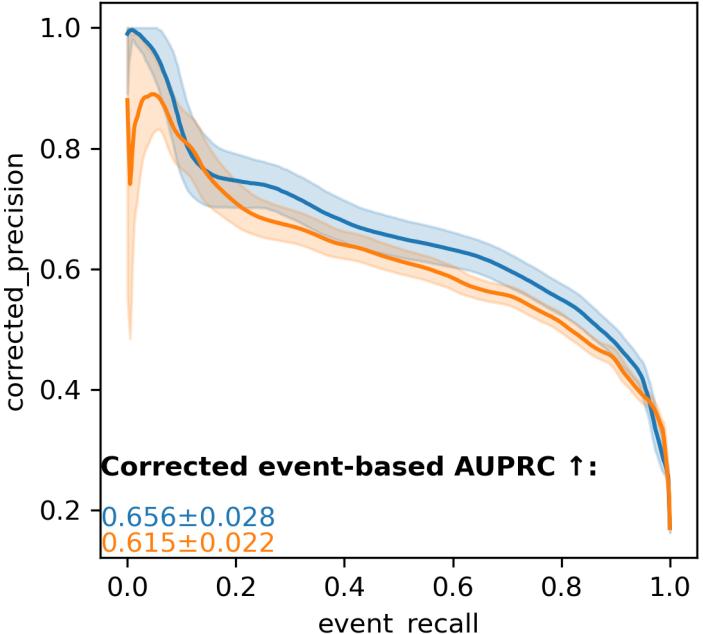
Corrected precision / recall curve



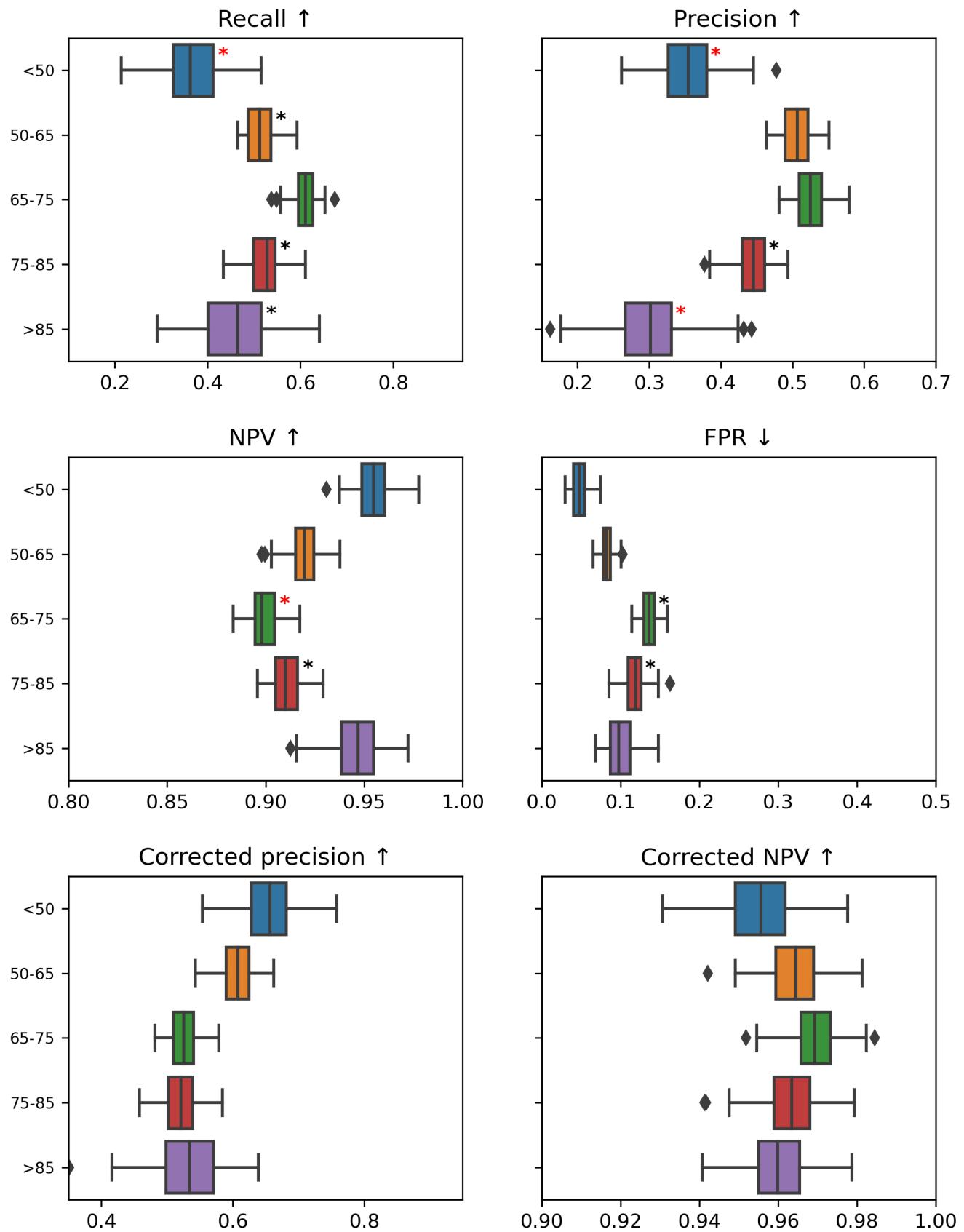
Precision / event-based recall curve

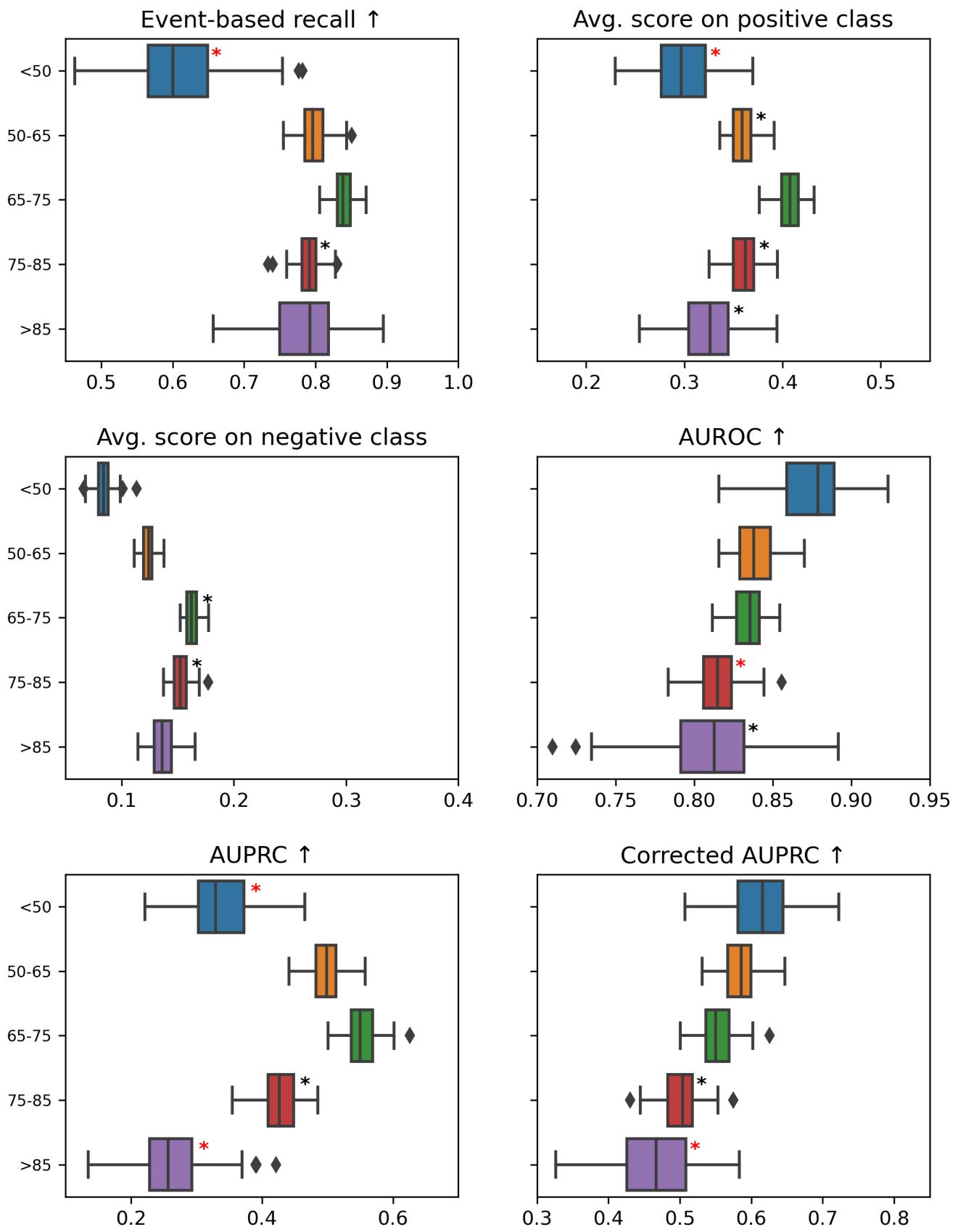


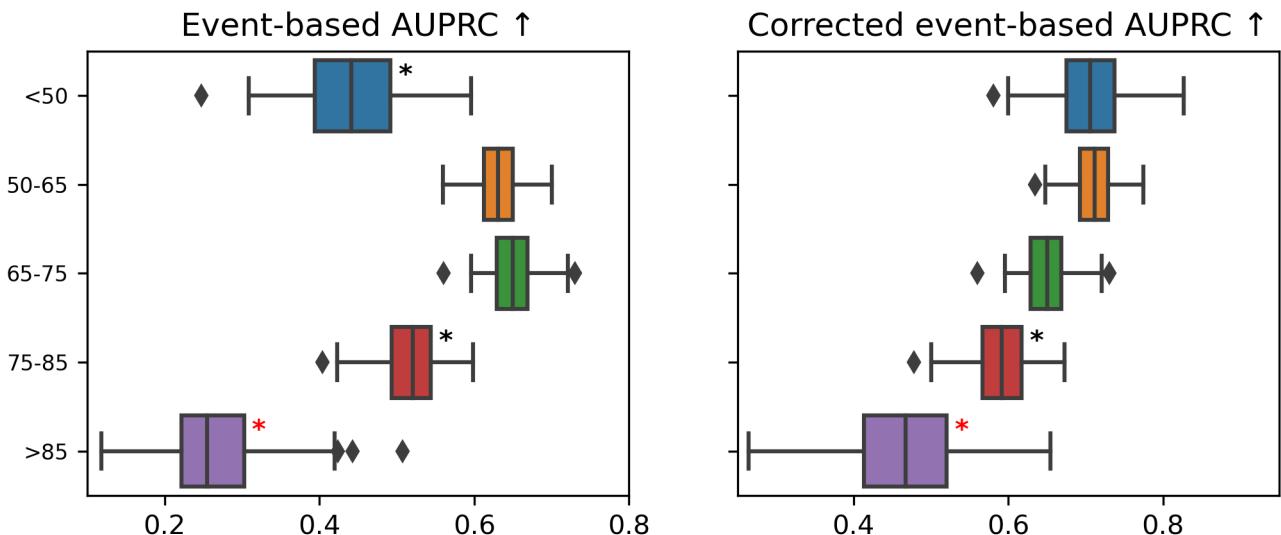
Corrected precision / event-based recall curve



Grouping by age_group



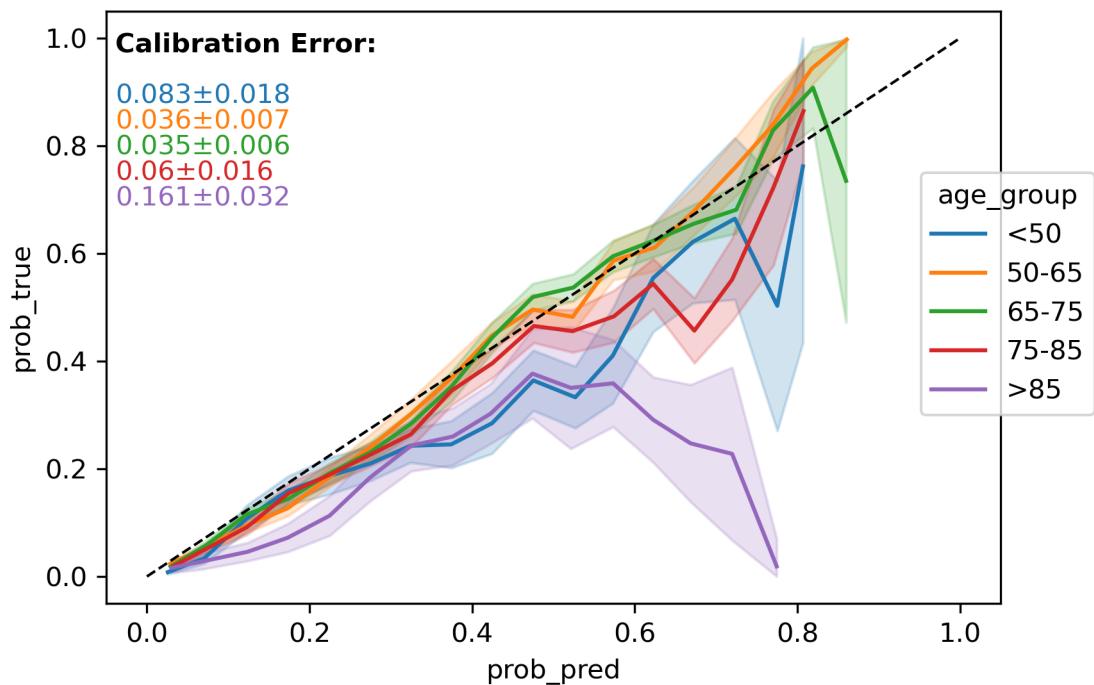




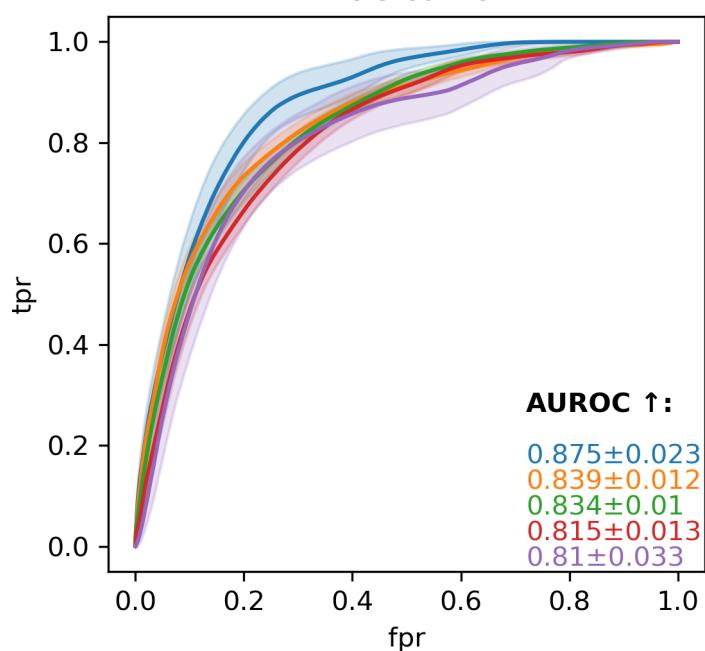
Metric	Group name	Group vs. rest	P-value	Delta
Recall ↑	<50	worse	1.36e-34	0.191
Recall ↑	50-65	worse	1.15e-19	0.043
Recall ↑	65-75	better	1.32e-34	0.111
Recall ↑	75-85	worse	4.78e-09	0.017
Recall ↑	>85	worse	1.34e-18	0.079
Precision ↑	<50	worse	1.78e-34	0.138
Precision ↑	50-65	better	1.56e-23	0.033
Precision ↑	65-75	better	1.36e-34	0.068
Precision ↑	75-85	worse	6.55e-32	0.048
Precision ↑	>85	worse	1.28e-34	0.188
NPV ↑	<50	better	1.28e-34	0.042
NPV ↑	65-75	worse	1.53e-34	0.028
NPV ↑	75-85	worse	9.19e-23	0.012
NPV ↑	>85	better	1.13e-31	0.029
FPR ↓	<50	better	1.28e-34	0.061
FPR ↓	50-65	better	1.03e-33	0.025
FPR ↓	65-75	worse	1.28e-34	0.051
FPR ↓	75-85	worse	4.82e-29	0.025
Corrected precision ↑	<50	better	1.28e-34	0.164
Corrected precision ↑	50-65	better	1.28e-34	0.113
Corrected NPV ↑	<50	better	9.19e-34	0.023
Corrected NPV ↑	50-65	better	1.28e-34	0.031
Corrected NPV ↑	65-75	better	1.28e-34	0.043
Corrected NPV ↑	75-85	better	1.28e-34	0.031
Corrected NPV ↑	>85	better	1.28e-34	0.028
Event-based recall ↑	<50	worse	1.28e-34	0.213
Event-based recall ↑	65-75	better	1.53e-34	0.063
Event-based recall ↑	75-85	worse	3.59e-08	0.011
Avg. score on positive class	<50	worse	1.78e-34	0.081
Avg. score on positive class	50-65	worse	2.76e-22	0.021

Avg. score on positive class	65-75	better	1.28e-34	0.055
Avg. score on positive class	75-85	worse	2.35e-18	0.014
Avg. score on positive class	>85	worse	4.66e-27	0.048
Avg. score on negative class	<50	better	1.28e-34	0.06
Avg. score on negative class	50-65	better	6.42e-33	0.016
Avg. score on negative class	65-75	worse	1.28e-34	0.039
Avg. score on negative class	75-85	worse	1.36e-34	0.023
AUROC ↑	<50	better	4.23e-28	0.045
AUROC ↑	75-85	worse	7.22e-33	0.034
AUROC ↑	>85	worse	1.40e-18	0.03
AUPRC ↑	<50	worse	1.32e-34	0.166
AUPRC ↑	50-65	better	6.94e-08	0.018
AUPRC ↑	65-75	better	1.28e-34	0.105
AUPRC ↑	75-85	worse	5.71e-34	0.074
AUPRC ↑	>85	worse	1.28e-34	0.236
Corrected AUPRC ↑	<50	better	2.71e-34	0.121
Corrected AUPRC ↑	50-65	better	1.95e-34	0.086
Corrected AUPRC ↑	65-75	better	5.53e-29	0.045
Corrected AUPRC ↑	75-85	worse	1.35e-14	0.025
Corrected AUPRC ↑	>85	worse	3.45e-10	0.046
Event-based AUPRC ↑	<50	worse	7.03e-34	0.165
Event-based AUPRC ↑	50-65	better	9.79e-24	0.048
Event-based AUPRC ↑	65-75	better	5.54e-34	0.09
Event-based AUPRC ↑	75-85	worse	8.66e-34	0.093
Event-based AUPRC ↑	>85	worse	1.28e-34	0.351
Corrected event-based AUPRC ↑	<50	better	2.11e-29	0.099
Corrected event-based AUPRC ↑	50-65	better	1.89e-34	0.108
Corrected event-based AUPRC ↑	65-75	better	2.55e-14	0.032
Corrected event-based AUPRC ↑	75-85	worse	3.26e-22	0.045
Corrected event-based AUPRC ↑	>85	worse	1.30e-31	0.157

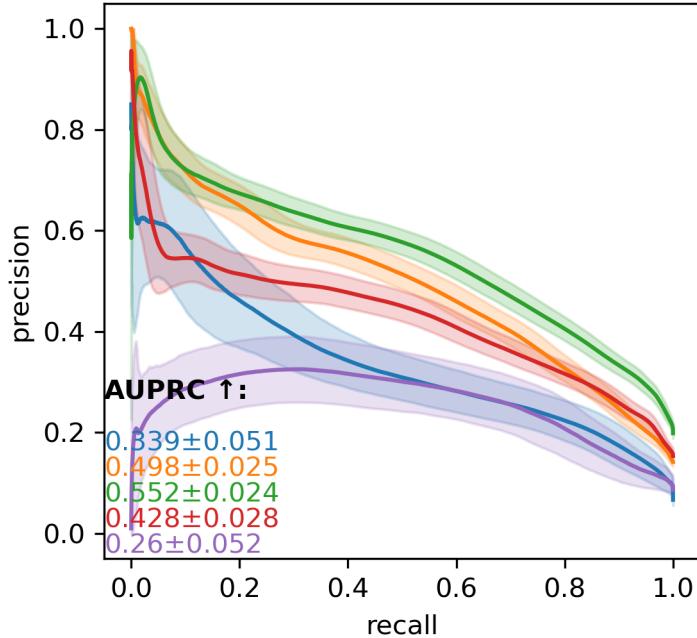
Calibration curve



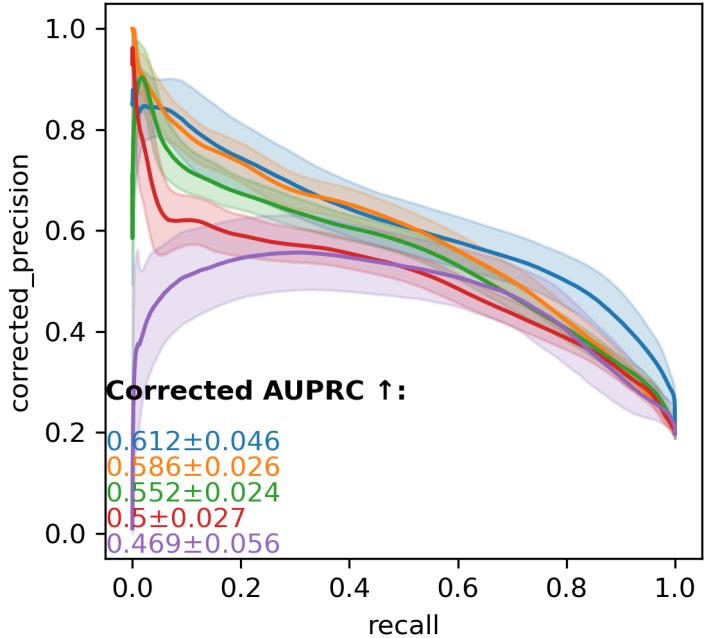
ROC curve



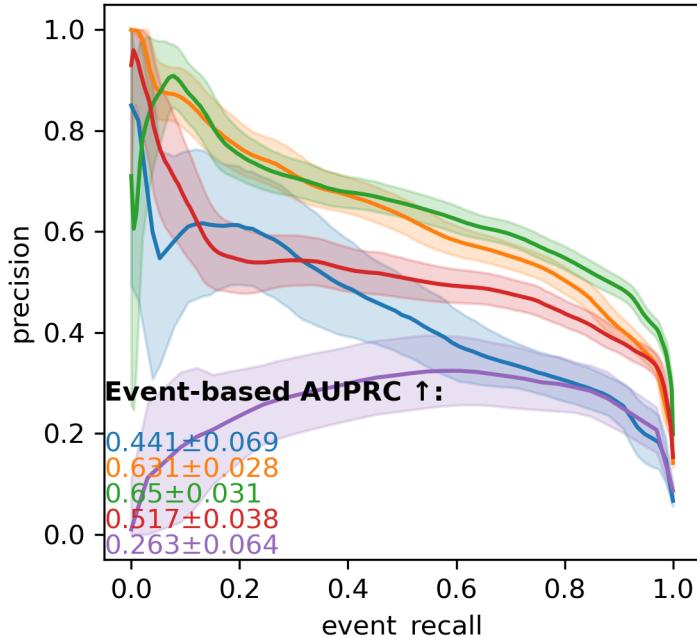
Precision / recall curve



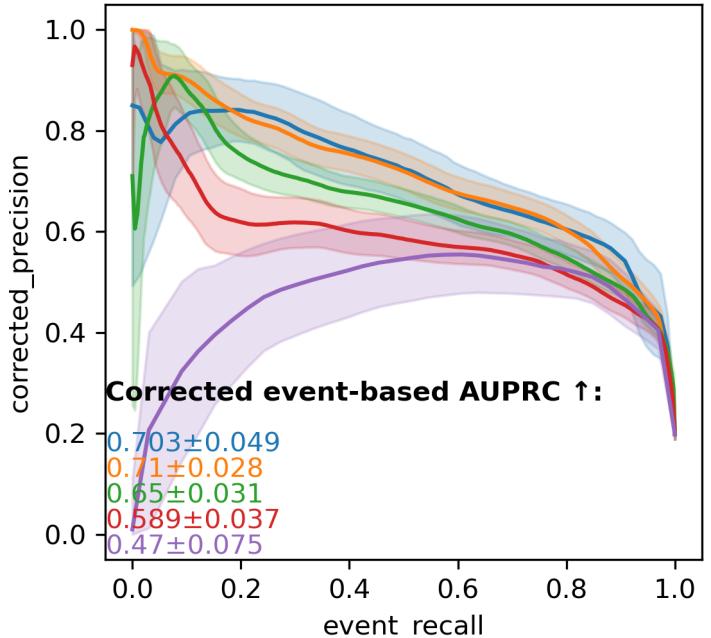
Corrected precision / recall curve



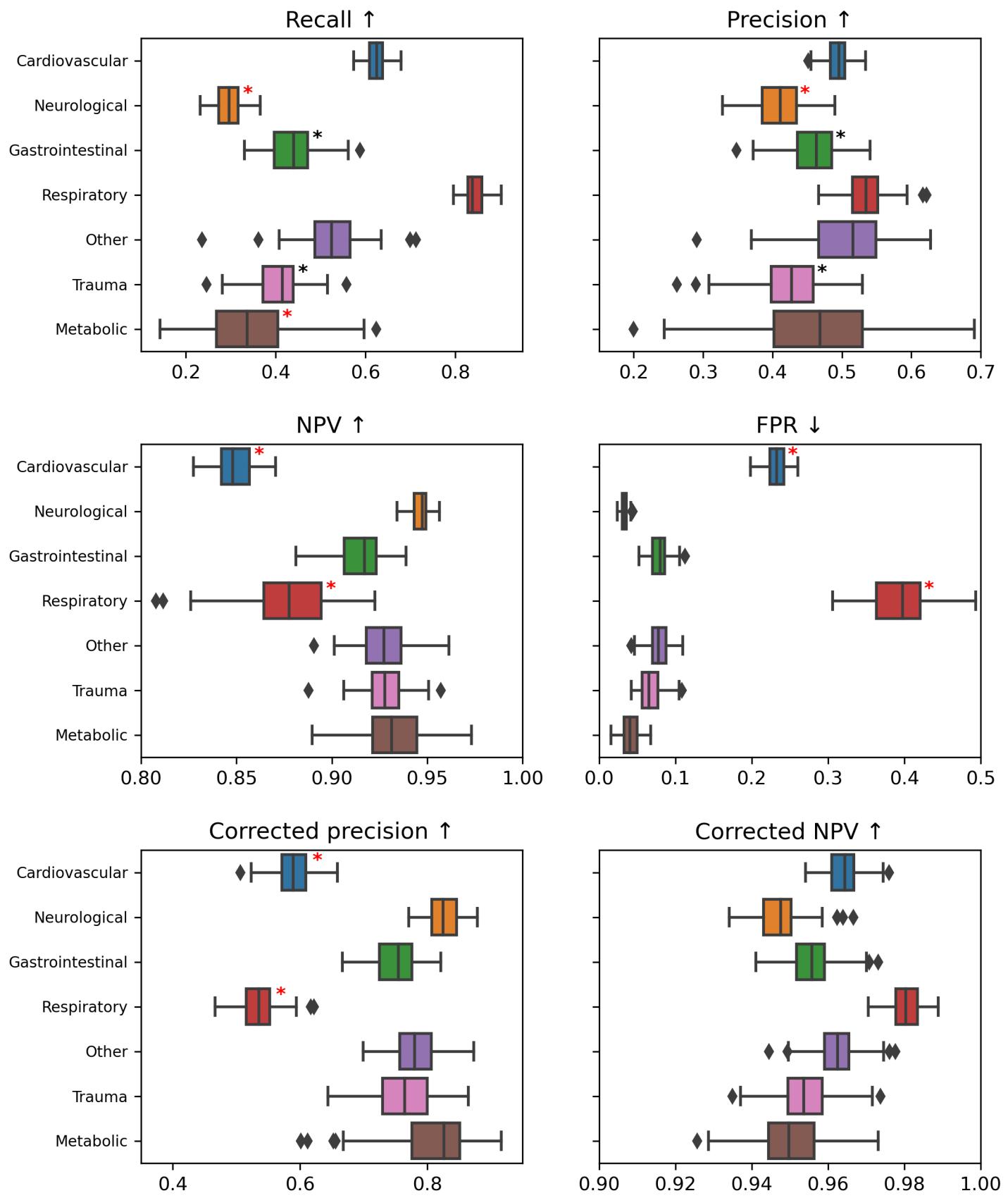
Precision / event-based recall curve



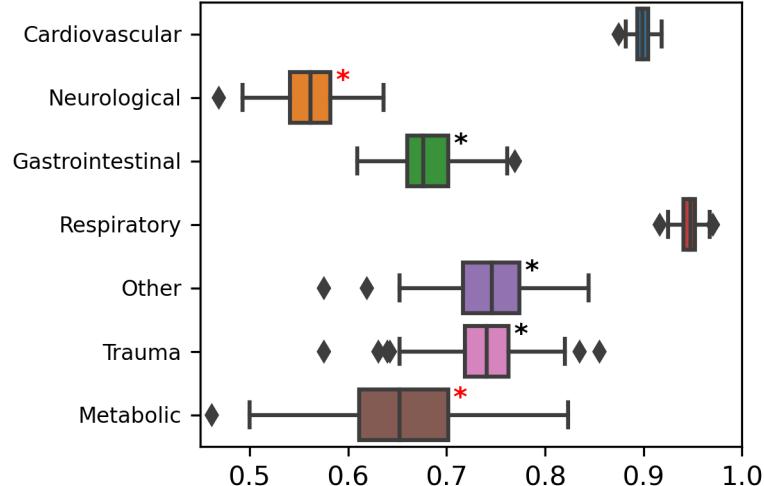
Corrected precision / event-based recall curve



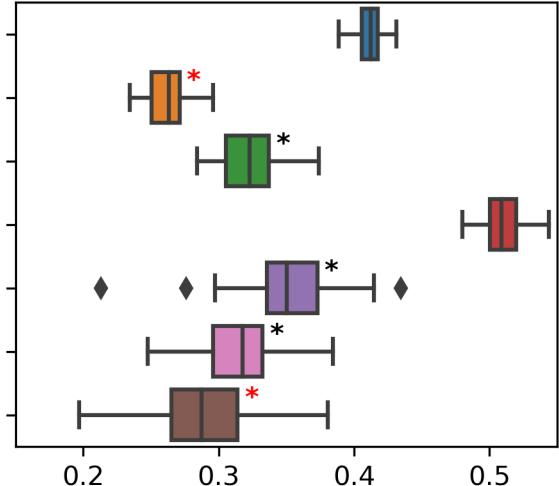
Grouping by APACHE_group



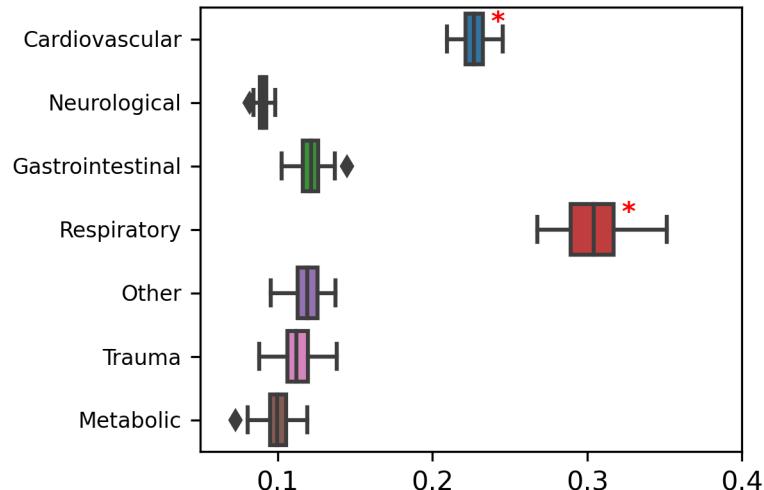
Event-based recall ↑



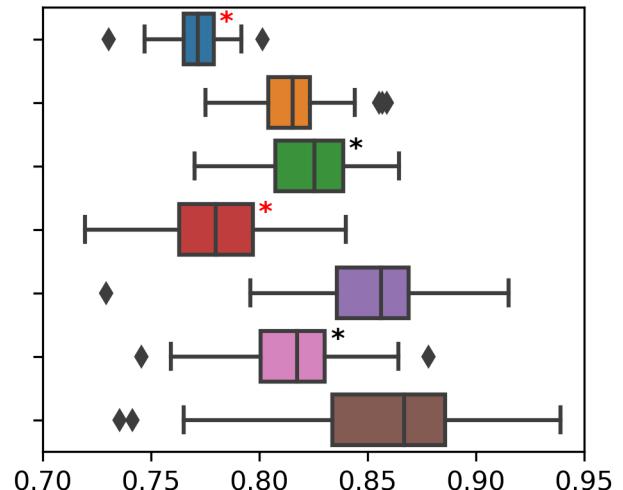
Avg. score on positive class



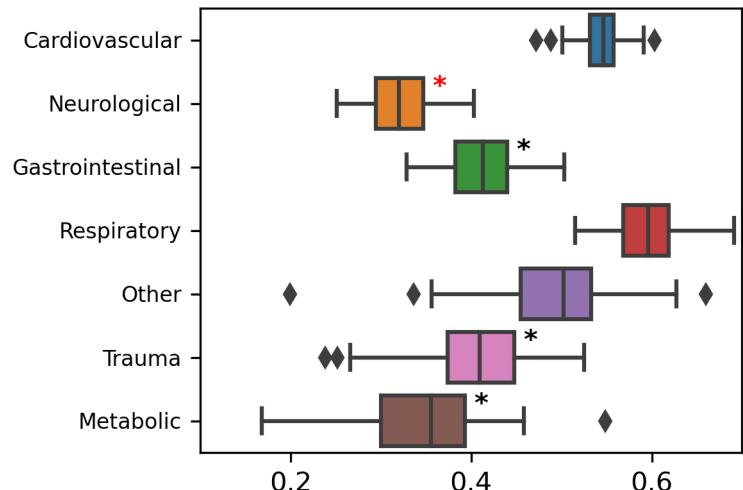
Avg. score on negative class



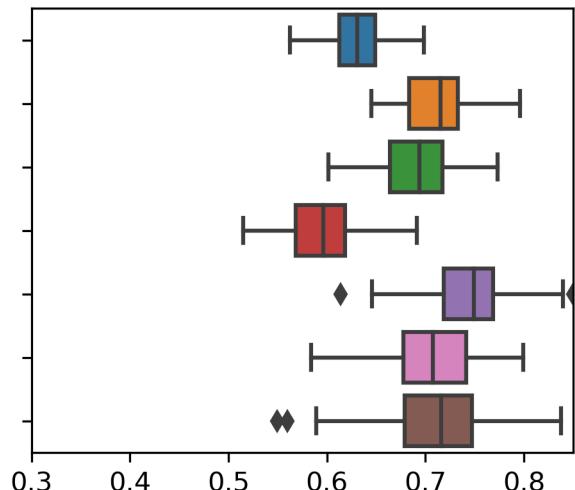
AUROC ↑



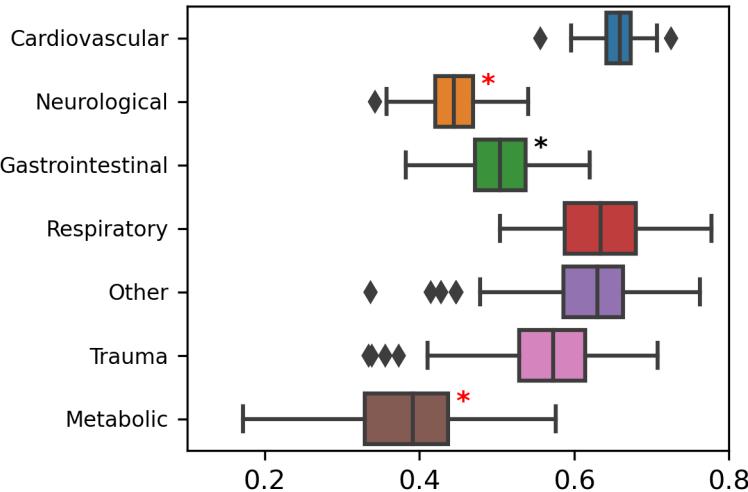
AUPRC ↑



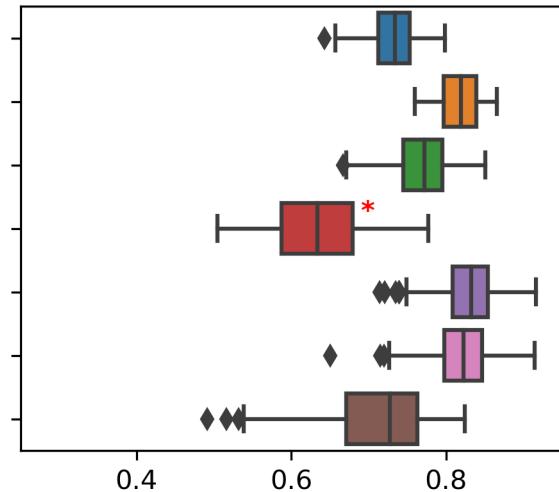
Corrected AUPRC ↑



Event-based AUPRC ↑



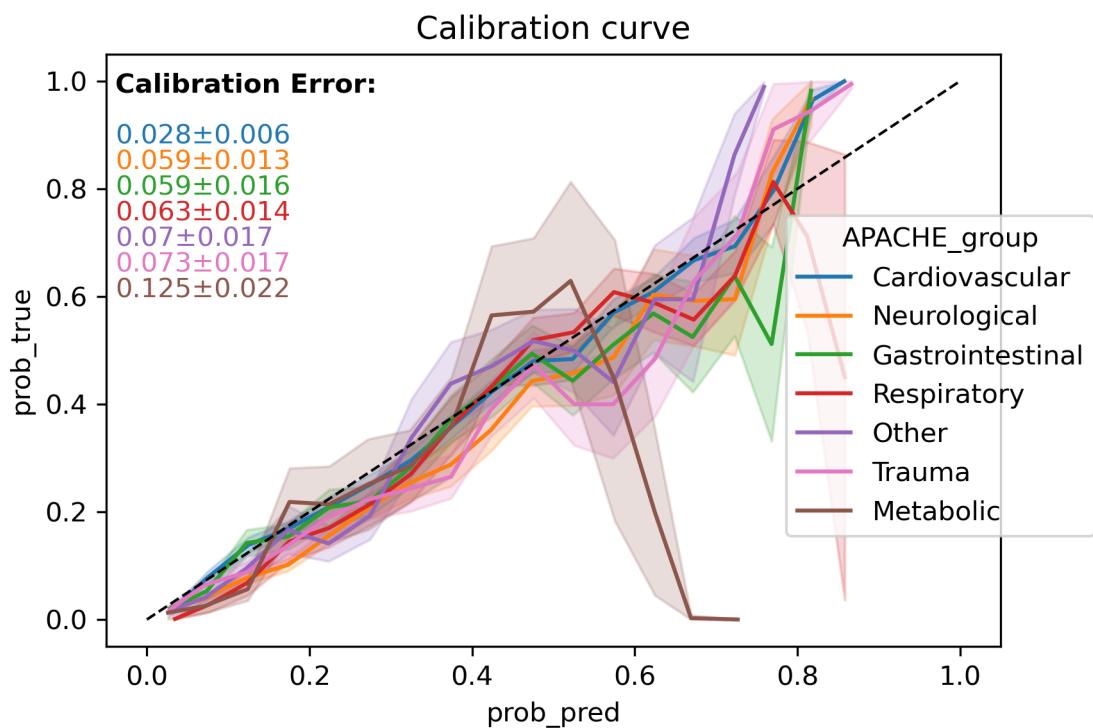
Corrected event-based AUPRC ↑

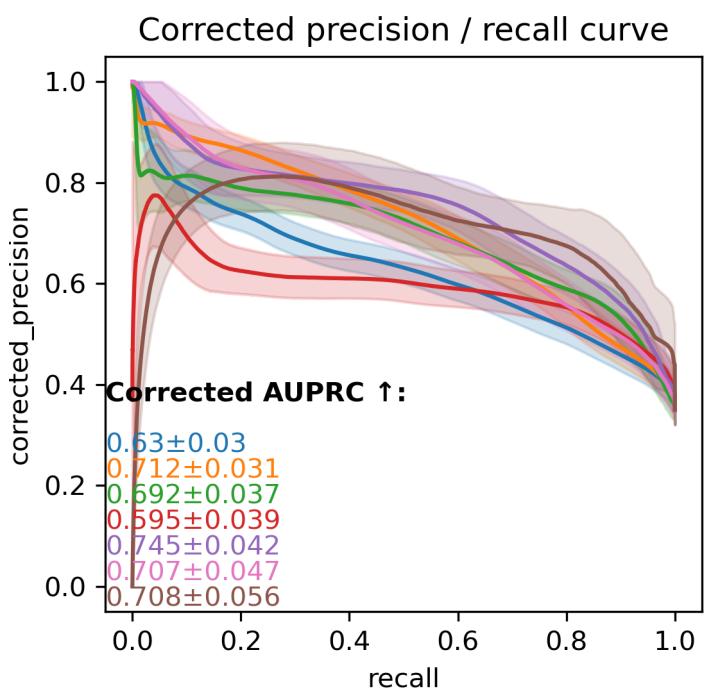
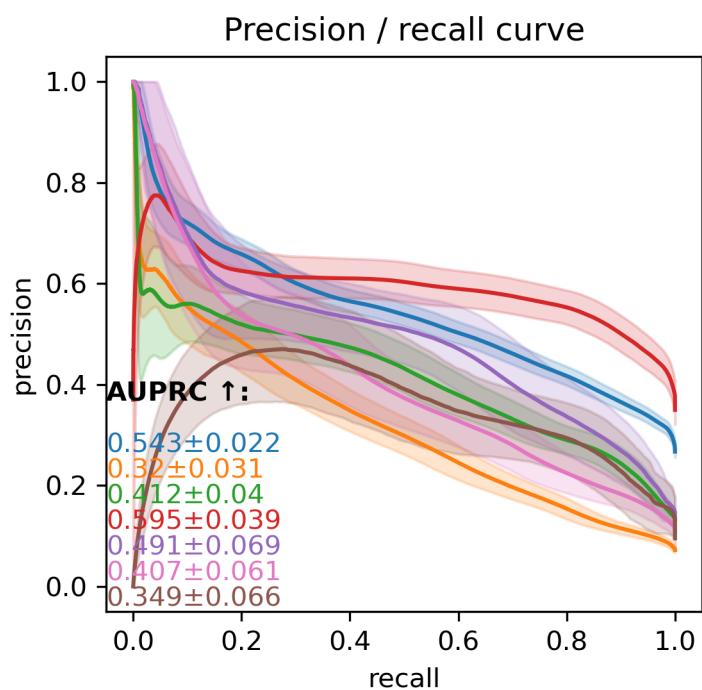
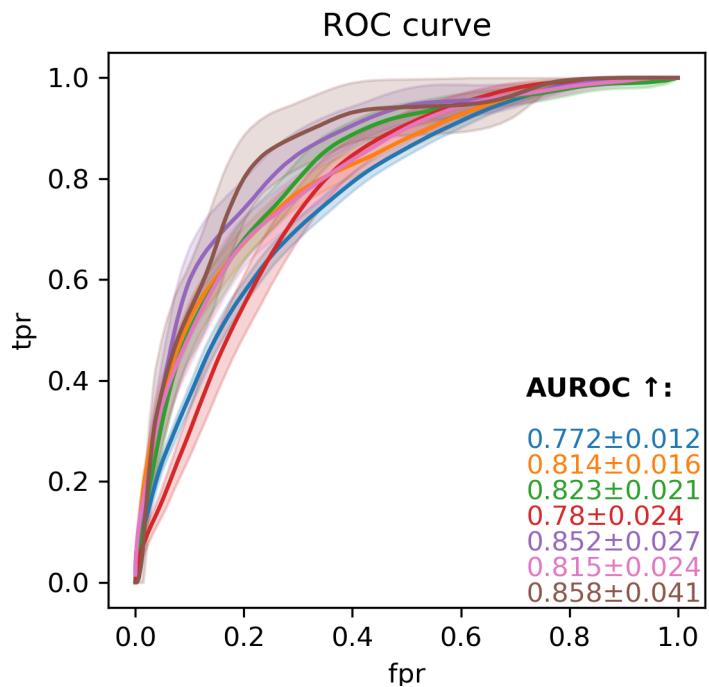


Metric	Group name	Group vs. rest	P-value	Delta
Recall ↑	Cardiovascular	better	1.28e-34	0.142
Recall ↑	Neurological	worse	1.28e-34	0.311
Recall ↑	Gastrointestinal	worse	8.06e-31	0.112
Recall ↑	Respiratory	better	1.28e-34	0.342
Recall ↑	Trauma	worse	7.92e-34	0.14
Recall ↑	Metabolic	worse	3.31e-27	0.209
Precision ↑	Cardiovascular	better	1.66e-14	0.024
Precision ↑	Neurological	worse	8.40e-34	0.085
Precision ↑	Gastrointestinal	worse	4.08e-09	0.022
Precision ↑	Respiratory	better	1.01e-31	0.062
Precision ↑	Other	better	2.86e-06	0.033
Precision ↑	Trauma	worse	2.25e-21	0.06
NPV ↑	Cardiovascular	worse	1.28e-34	0.087
NPV ↑	Neurological	better	1.28e-34	0.054
NPV ↑	Respiratory	worse	5.88e-33	0.043
NPV ↑	Other	better	1.47e-07	0.008
NPV ↑	Trauma	better	2.15e-13	0.01
NPV ↑	Metabolic	better	2.45e-12	0.012
FPR ↓	Cardiovascular	worse	1.28e-34	0.166
FPR ↓	Neurological	better	1.28e-34	0.123
FPR ↓	Gastrointestinal	better	2.49e-29	0.022
FPR ↓	Respiratory	worse	1.28e-34	0.312
FPR ↓	Other	better	8.06e-28	0.023
FPR ↓	Trauma	better	6.00e-32	0.038
FPR ↓	Metabolic	better	1.28e-34	0.06
Corrected precision ↑	Cardiovascular	worse	1.75e-27	0.058
Corrected precision ↑	Neurological	better	1.28e-34	0.329
Corrected precision ↑	Gastrointestinal	better	1.28e-34	0.174
Corrected precision ↑	Respiratory	worse	1.19e-28	0.06
Corrected precision ↑	Other	better	1.28e-34	0.202

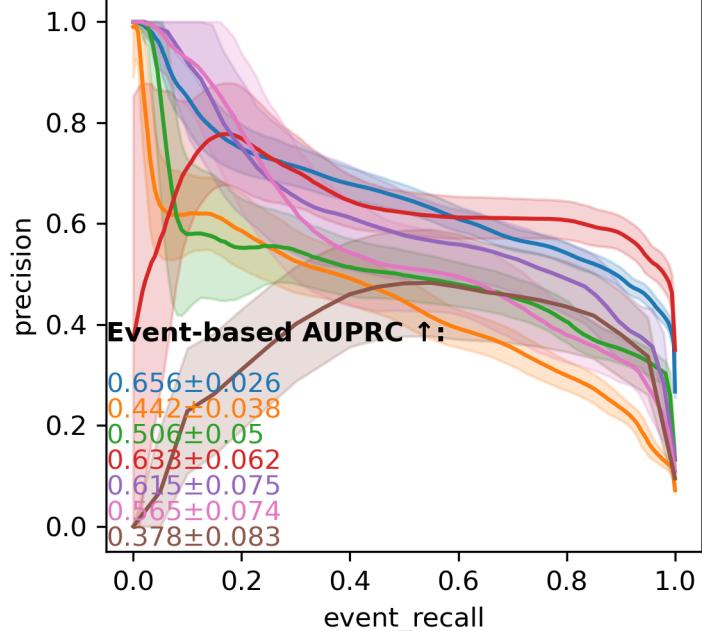
Corrected precision ↑	Trauma	better	1.28e-34	0.186
Corrected precision ↑	Metabolic	better	1.49e-34	0.247
Corrected NPV ↑	Cardiovascular	better	1.28e-34	0.029
Corrected NPV ↑	Gastrointestinal	better	2.83e-32	0.014
Corrected NPV ↑	Respiratory	better	1.28e-34	0.045
Corrected NPV ↑	Other	better	1.78e-34	0.022
Corrected NPV ↑	Trauma	better	4.59e-28	0.012
Corrected NPV ↑	Metabolic	better	7.68e-16	0.009
Event-based recall ↑	Cardiovascular	better	1.28e-34	0.175
Event-based recall ↑	Neurological	worse	1.28e-34	0.296
Event-based recall ↑	Gastrointestinal	worse	1.28e-34	0.135
Event-based recall ↑	Respiratory	better	1.28e-34	0.172
Event-based recall ↑	Other	worse	3.77e-26	0.057
Event-based recall ↑	Trauma	worse	1.89e-27	0.065
Event-based recall ↑	Metabolic	worse	3.19e-32	0.15
Avg. score on positive class	Cardiovascular	better	1.28e-34	0.067
Avg. score on positive class	Neurological	worse	1.28e-34	0.139
Avg. score on positive class	Gastrointestinal	worse	4.92e-34	0.055
Avg. score on positive class	Respiratory	better	1.28e-34	0.157
Avg. score on positive class	Other	worse	3.4e-11	0.024
Avg. score on positive class	Trauma	worse	7.01e-33	0.06
Avg. score on positive class	Metabolic	worse	2.59e-31	0.087
Avg. score on negative class	Cardiovascular	worse	1.28e-34	0.115
Avg. score on negative class	Neurological	better	1.28e-34	0.081
Avg. score on negative class	Gastrointestinal	better	1.15e-29	0.014
Avg. score on negative class	Respiratory	worse	1.28e-34	0.178
Avg. score on negative class	Other	better	1.69e-28	0.016
Avg. score on negative class	Trauma	better	3.44e-31	0.025
Avg. score on negative class	Metabolic	better	1.28e-34	0.035
AUROC ↑	Cardiovascular	worse	1.28e-34	0.074
AUROC ↑	Gastrointestinal	worse	9.15e-14	0.019
AUROC ↑	Respiratory	worse	9.75e-34	0.057
AUROC ↑	Other	better	1.01e-06	0.015
AUROC ↑	Trauma	worse	2.08e-22	0.026
AUROC ↑	Metabolic	better	1.54e-06	0.025
AUPRC ↑	Cardiovascular	better	1.44e-34	0.107
AUPRC ↑	Neurological	worse	1.28e-34	0.203
AUPRC ↑	Gastrointestinal	worse	2.31e-31	0.079
AUPRC ↑	Respiratory	better	1.28e-34	0.133
AUPRC ↑	Trauma	worse	6.93e-24	0.082
AUPRC ↑	Metabolic	worse	2.91e-33	0.131
Corrected AUPRC ↑	Cardiovascular	better	2.97e-15	0.035
Corrected AUPRC ↑	Neurological	better	1.28e-34	0.193

Corrected AUPRC ↑	Gastrointestinal	better	1.68e-34	0.117
Corrected AUPRC ↑	Respiratory	better	3.05e-07	0.022
Corrected AUPRC ↑	Other	better	1.28e-34	0.176
Corrected AUPRC ↑	Trauma	better	4.78e-34	0.134
Corrected AUPRC ↑	Metabolic	better	6.93e-32	0.142
Event-based AUPRC ↑	Cardiovascular	better	5.38e-34	0.114
Event-based AUPRC ↑	Neurological	worse	1.28e-34	0.179
Event-based AUPRC ↑	Gastrointestinal	worse	4.44e-31	0.099
Event-based AUPRC ↑	Respiratory	better	1.91e-09	0.046
Event-based AUPRC ↑	Other	better	7.02e-07	0.035
Event-based AUPRC ↑	Metabolic	worse	1.73e-34	0.207
Corrected event-based AUPRC ↑	Cardiovascular	better	7.07e-14	0.037
Corrected event-based AUPRC ↑	Neurological	better	1.28e-34	0.196
Corrected event-based AUPRC ↑	Gastrointestinal	better	3.11e-29	0.087
Corrected event-based AUPRC ↑	Respiratory	worse	9.08e-13	0.058
Corrected event-based AUPRC ↑	Other	better	1.28e-34	0.154
Corrected event-based AUPRC ↑	Trauma	better	1.71e-33	0.15
Corrected event-based AUPRC ↑	Metabolic	better	1.95e-07	0.047

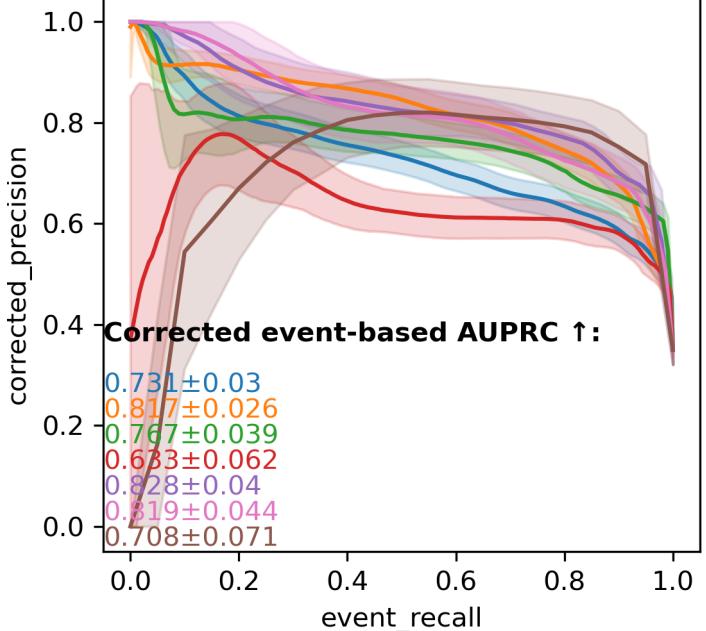




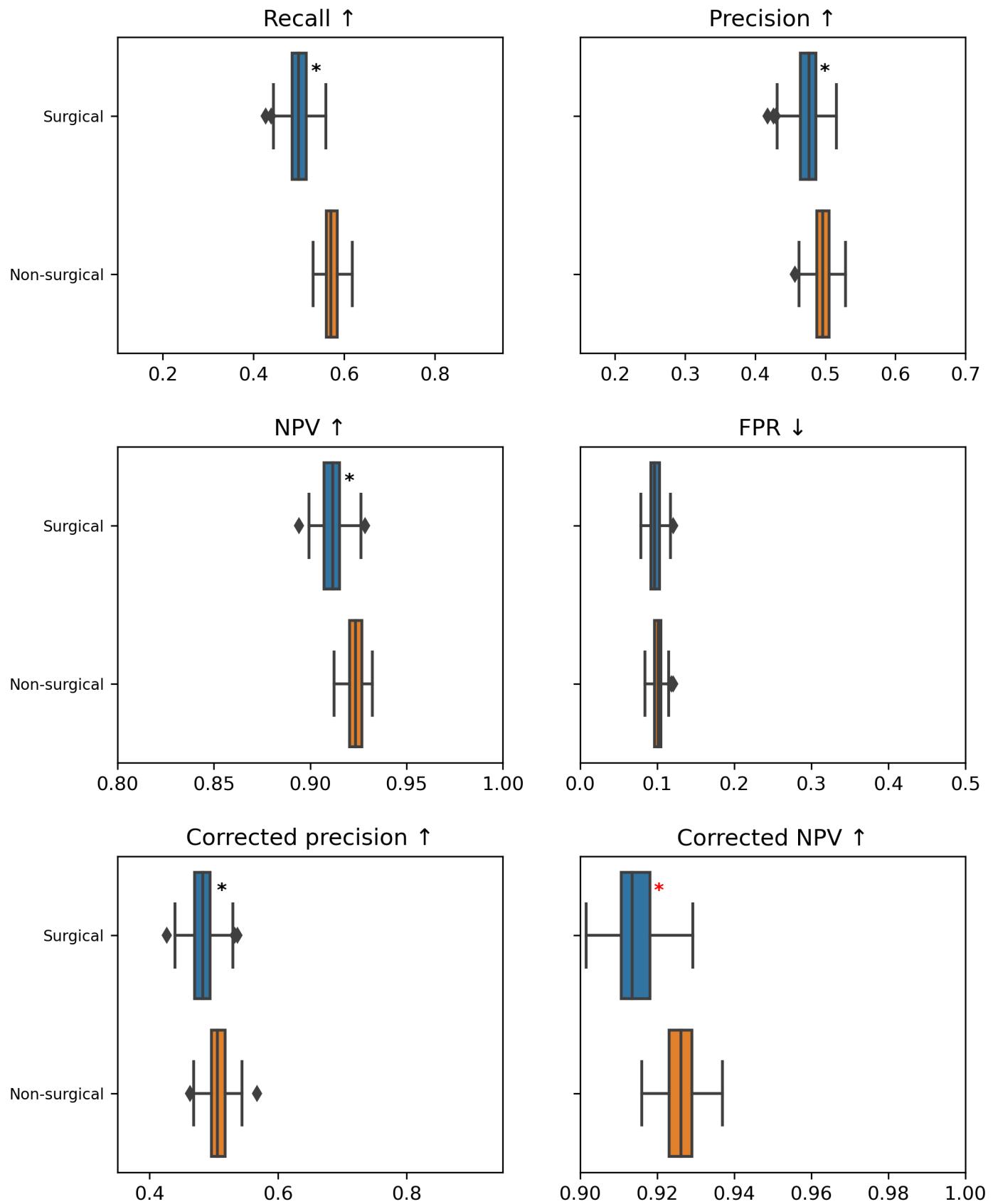
Precision / event-based recall curve



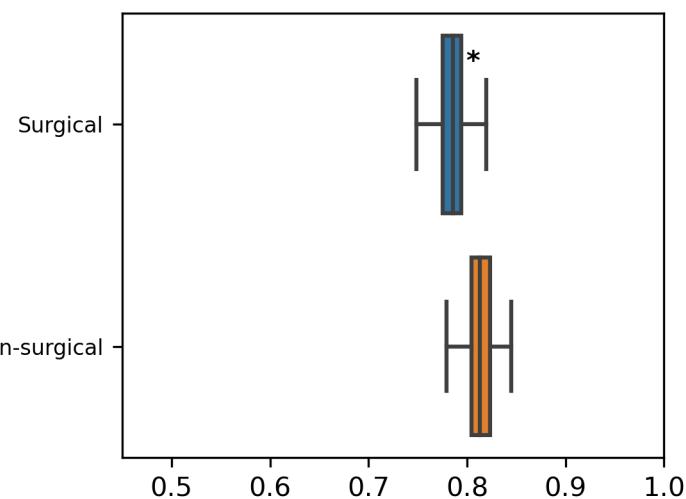
Corrected precision / event-based recall curve



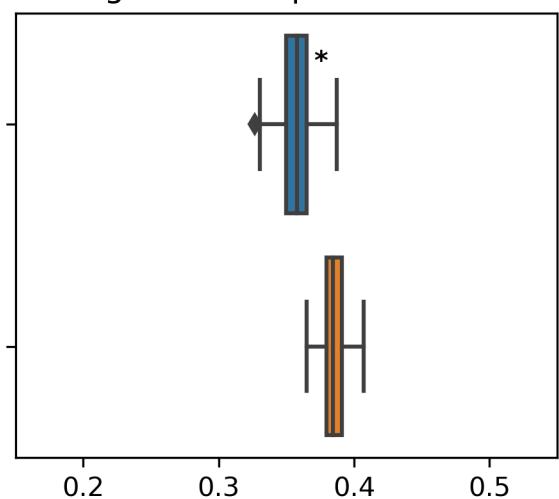
Grouping by surgical_status



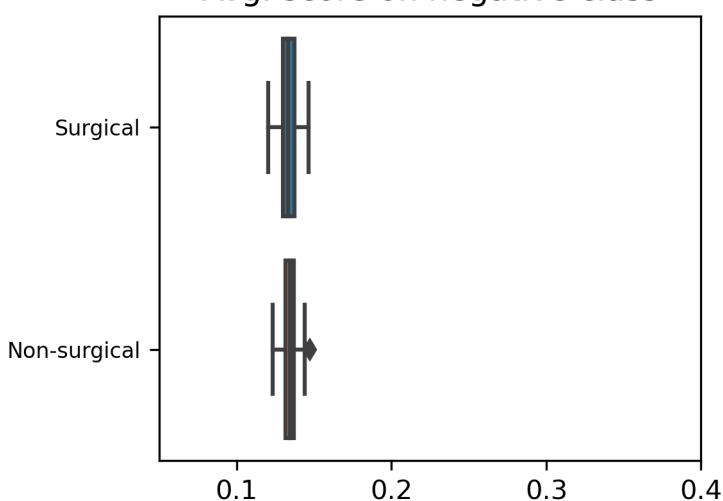
Event-based recall ↑



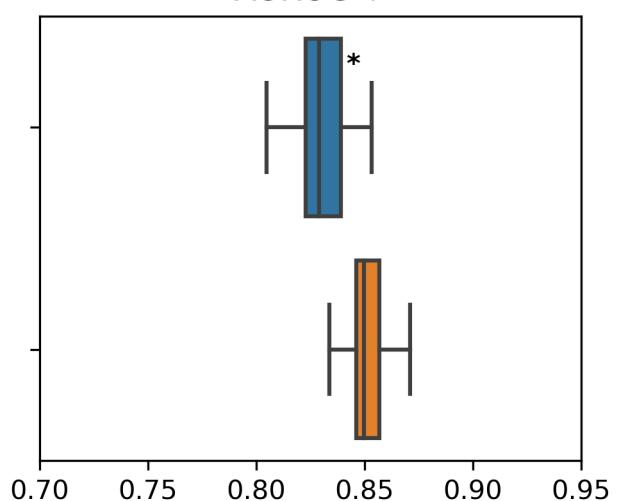
Avg. score on positive class



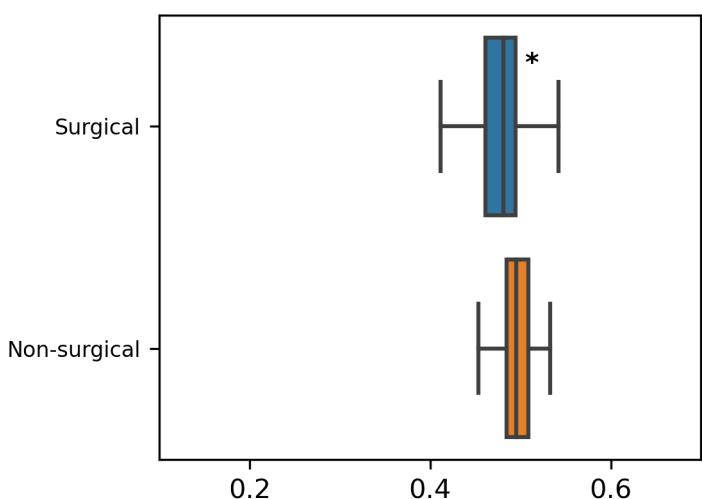
Avg. score on negative class



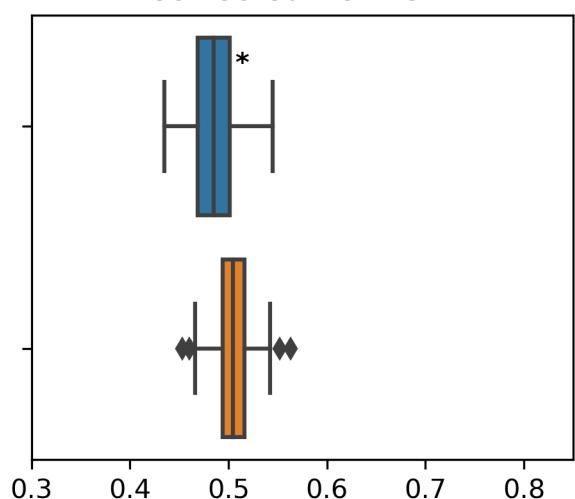
AUROC ↑

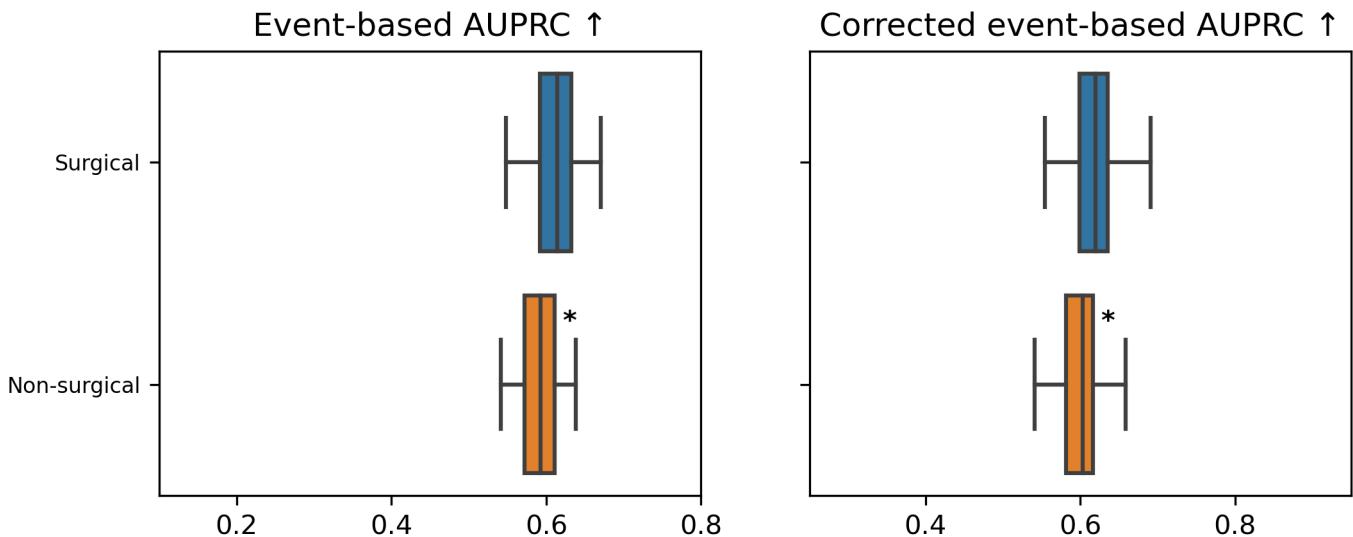


AUPRC ↑

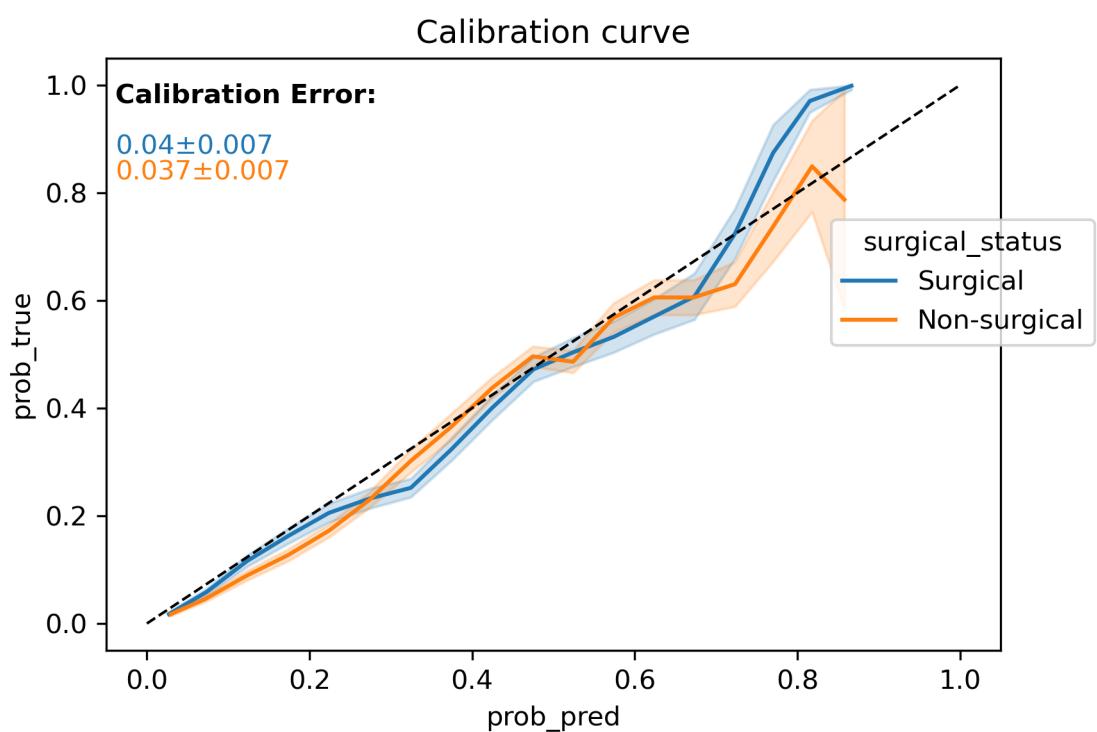


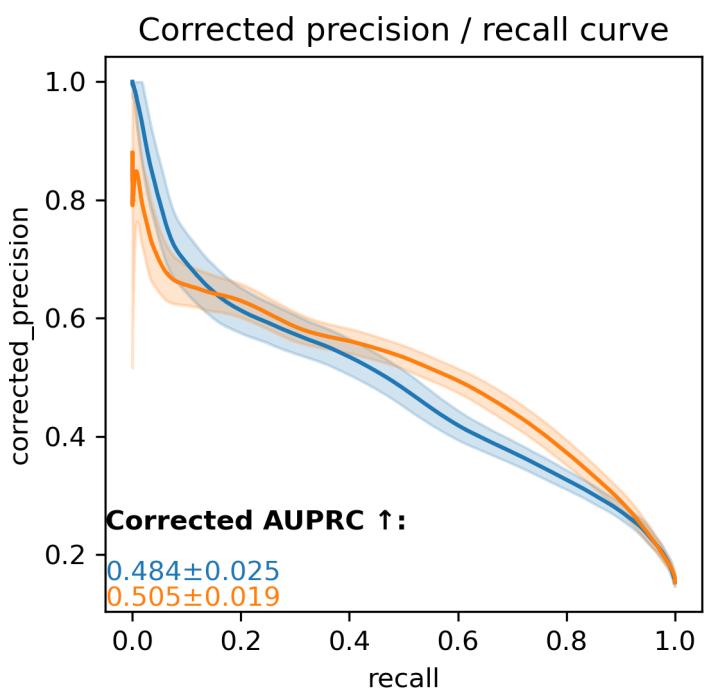
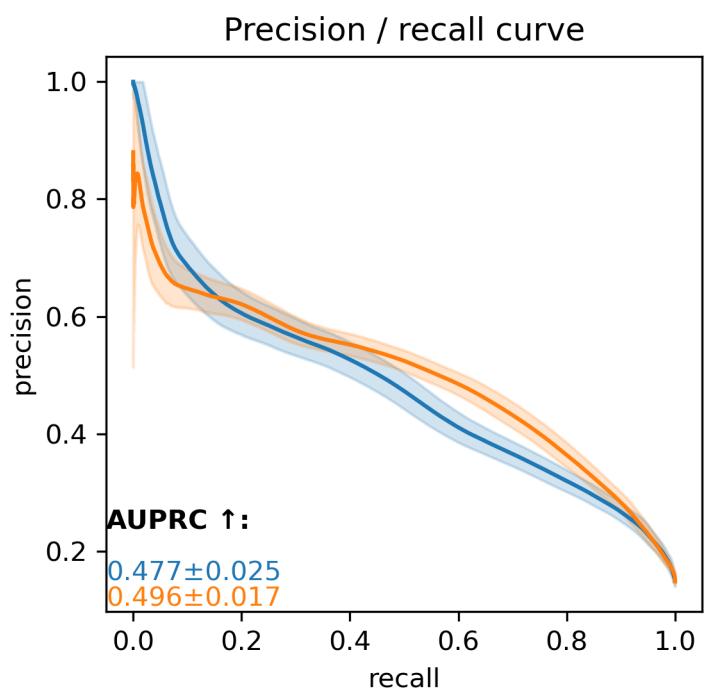
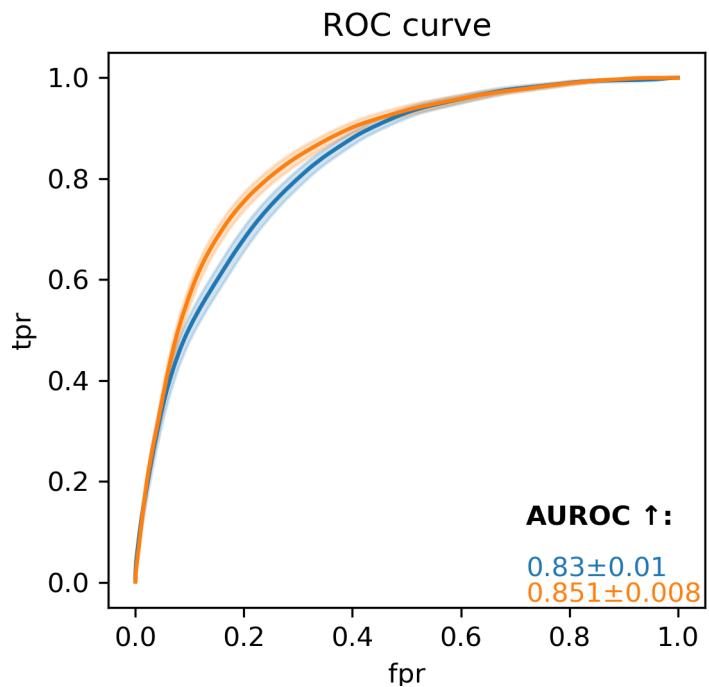
Corrected AUPRC ↑



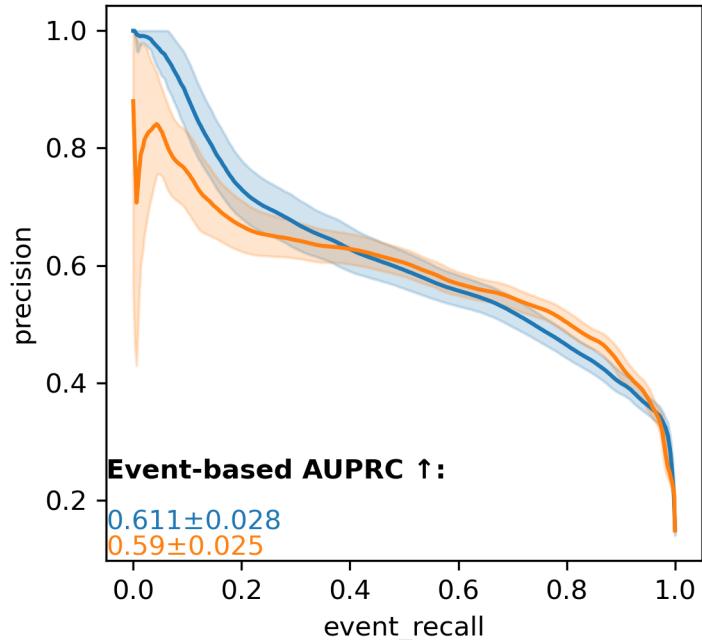


Metric	Group with worst dist.	P-value	Delta
Recall ↑	Surgical	2.17e-33	0.071
Precision ↑	Surgical	1.71e-16	0.019
NPV ↑	Surgical	4.89e-26	0.012
Corrected precision ↑	Surgical	2.72e-16	0.023
Corrected NPV ↑	Surgical	3.78e-26	0.013
Event-based recall ↑	Surgical	3.87e-24	0.027
Avg. score on positive class	Surgical	1.34e-31	0.027
AUROC ↑	Surgical	5.39e-28	0.021
AUPRC ↑	Surgical	4.17e-08	0.014
Corrected AUPRC ↑	Surgical	7.42e-10	0.02
Event-based AUPRC ↑	Non-surgical	2.28e-07	0.021
Corrected event-based AUPRC ↑	Non-surgical	1.88e-06	0.016

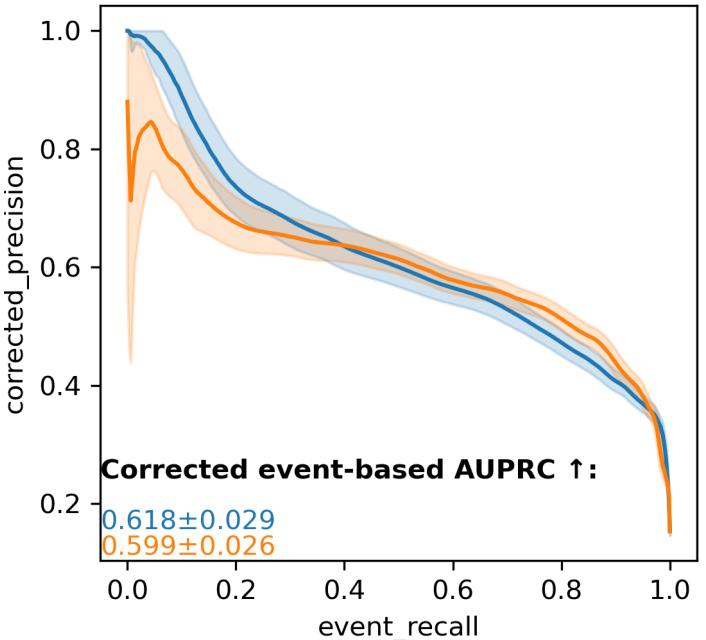




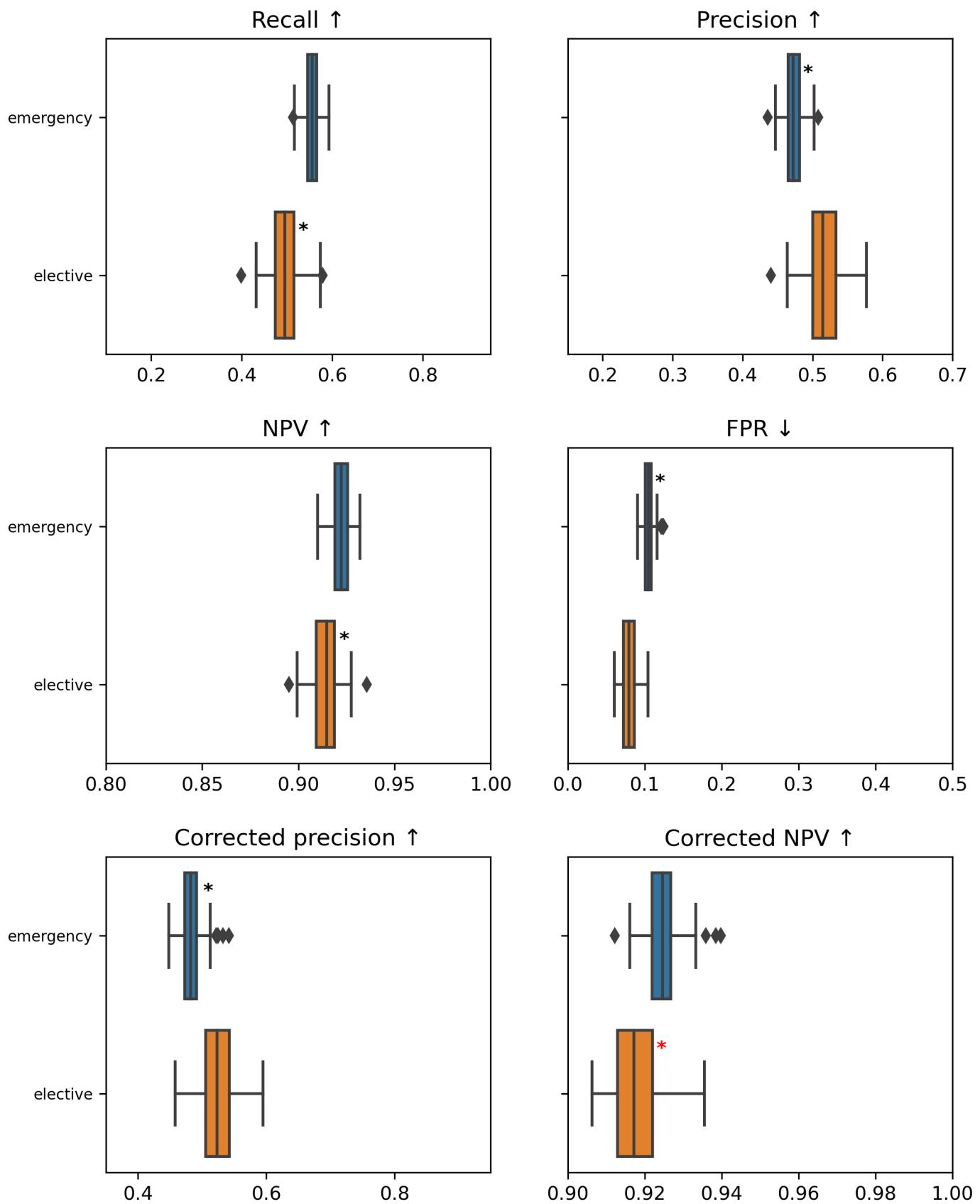
Precision / event-based recall curve



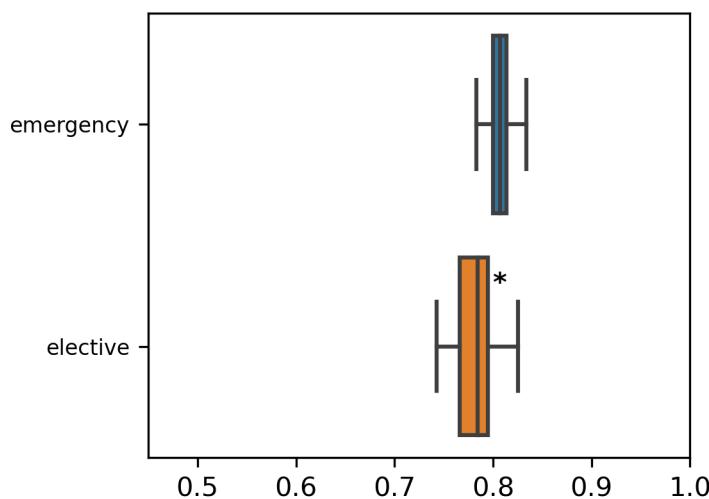
Corrected precision / event-based recall curve



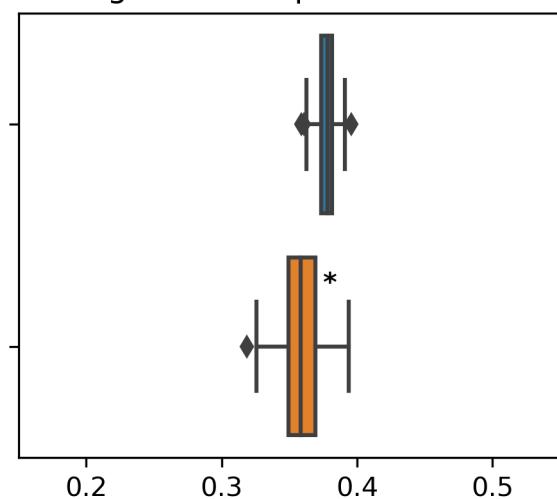
Grouping by admission_type



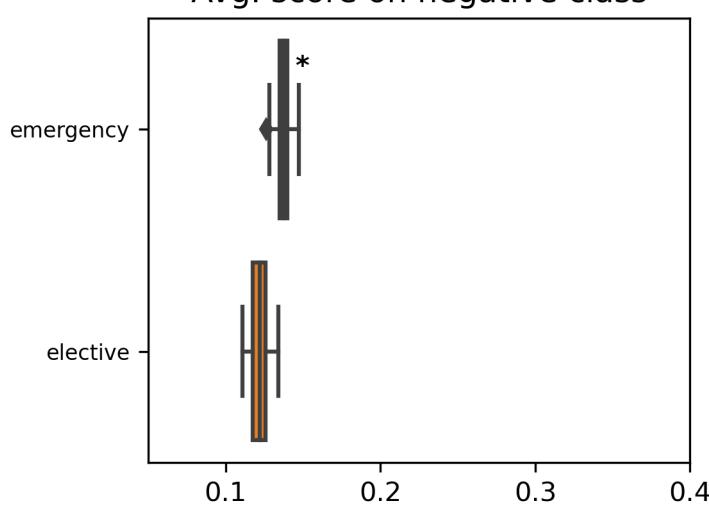
Event-based recall ↑



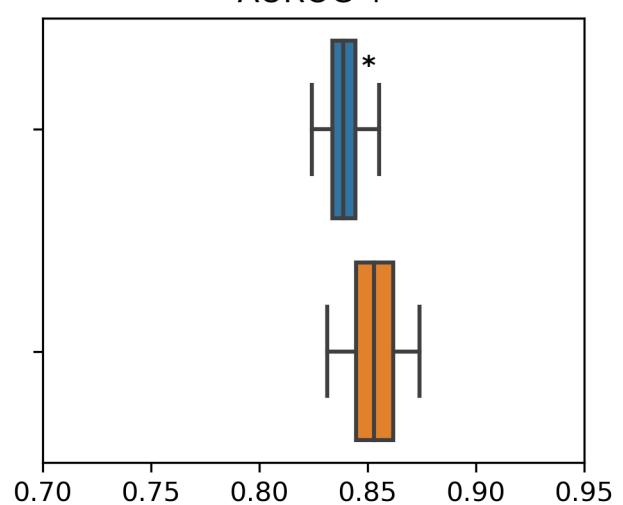
Avg. score on positive class



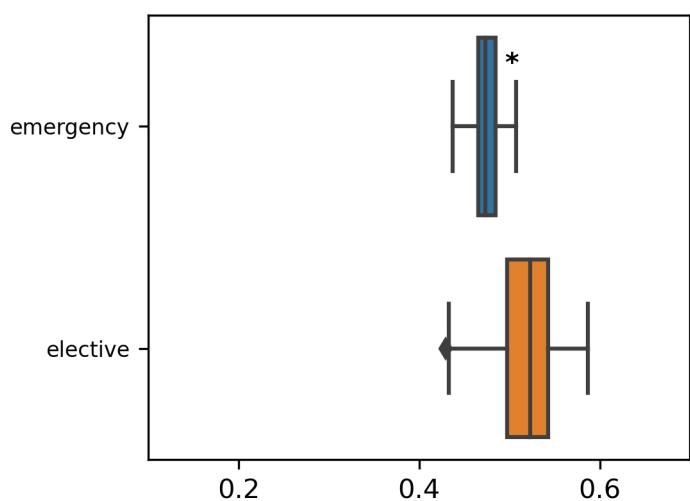
Avg. score on negative class



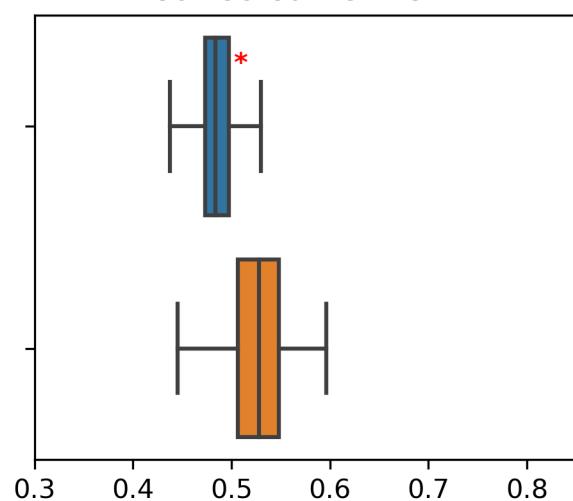
AUROC ↑

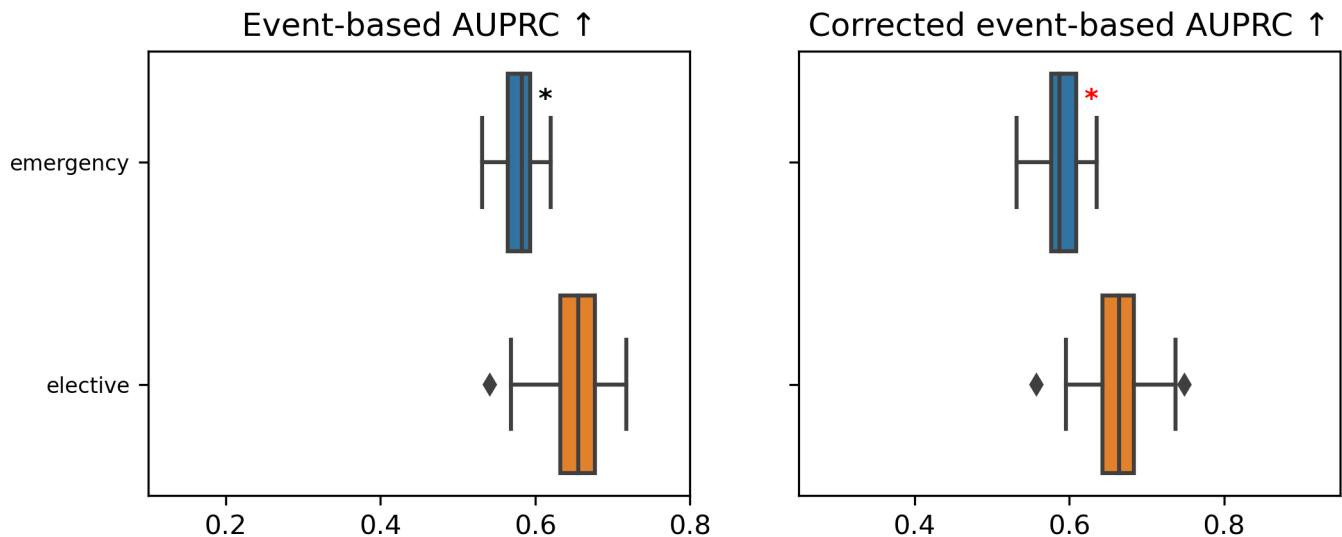


AUPRC ↑

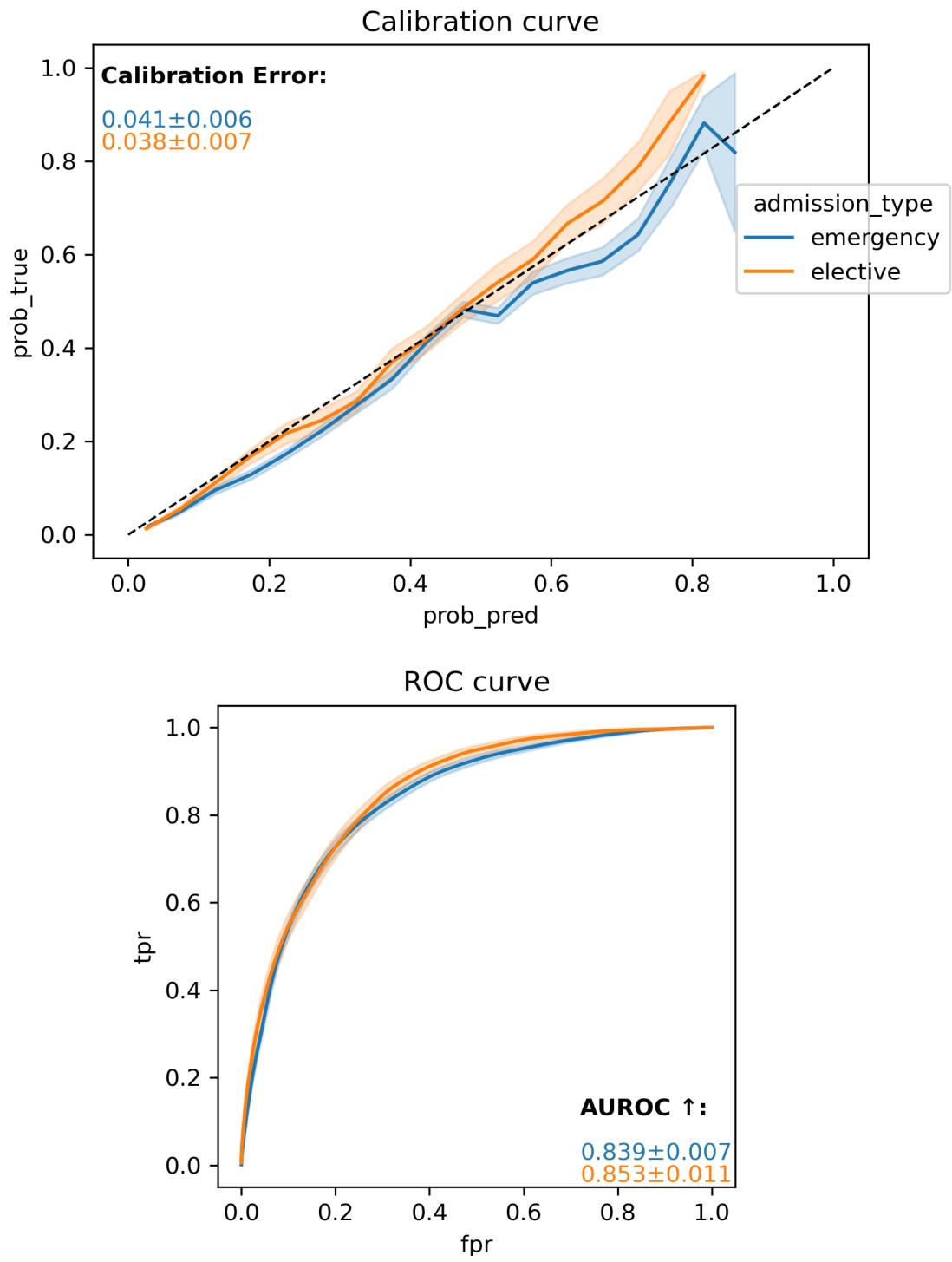


Corrected AUPRC ↑

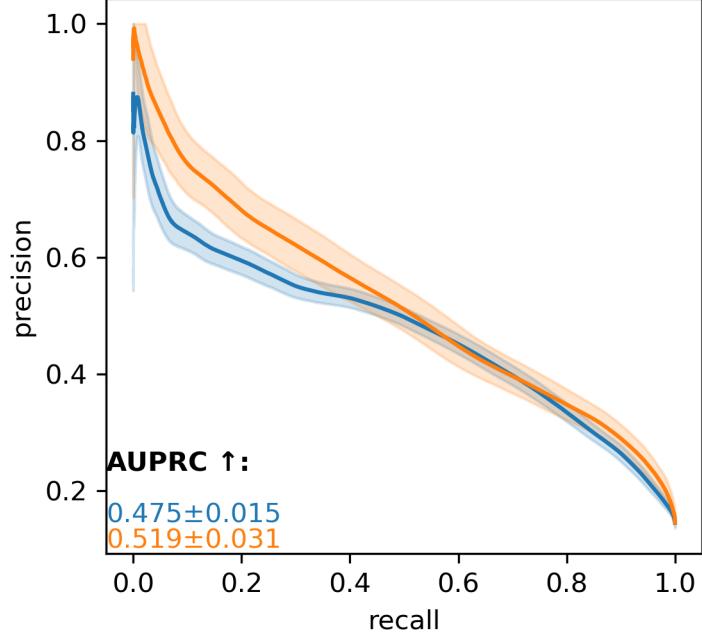




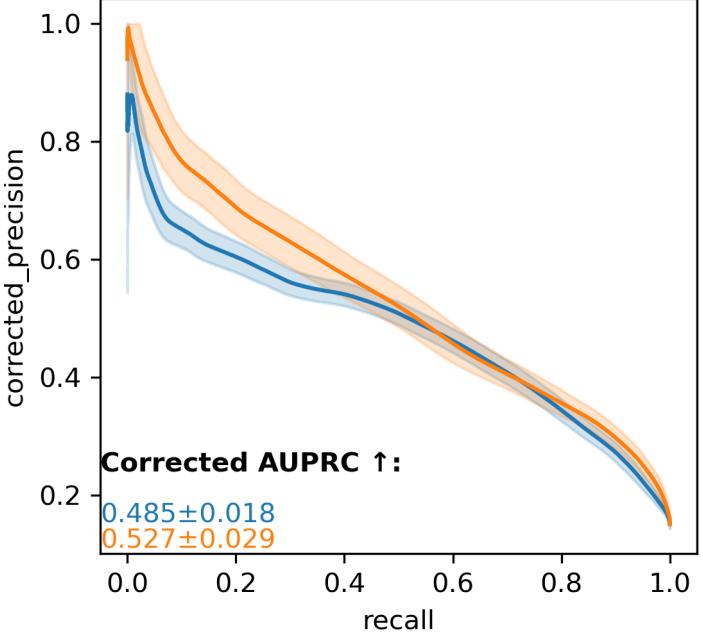
Metric	Group with worst dist.	P-value	Delta
Recall ↑	elective	7.05e-28	0.061
Precision ↑	emergency	5.71e-26	0.043
NPV ↑	elective	9.93e-16	0.008
FPR ↓	emergency	3.27e-33	0.025
Corrected precision ↑	emergency	2.63e-24	0.041
Corrected NPV ↑	elective	6.18e-16	0.007
Event-based recall ↑	elective	1.19e-20	0.023
Avg. score on positive class	elective	1.47e-19	0.019
Avg. score on negative class	emergency	1.87e-33	0.016
AUROC ↑	emergency	1.02e-17	0.014
AUPRC ↑	emergency	3.95e-22	0.049
Corrected AUPRC ↑	emergency	9.23e-22	0.044
Event-based AUPRC ↑	emergency	1.42e-30	0.073
Corrected event-based AUPRC ↑	emergency	3.89e-30	0.077



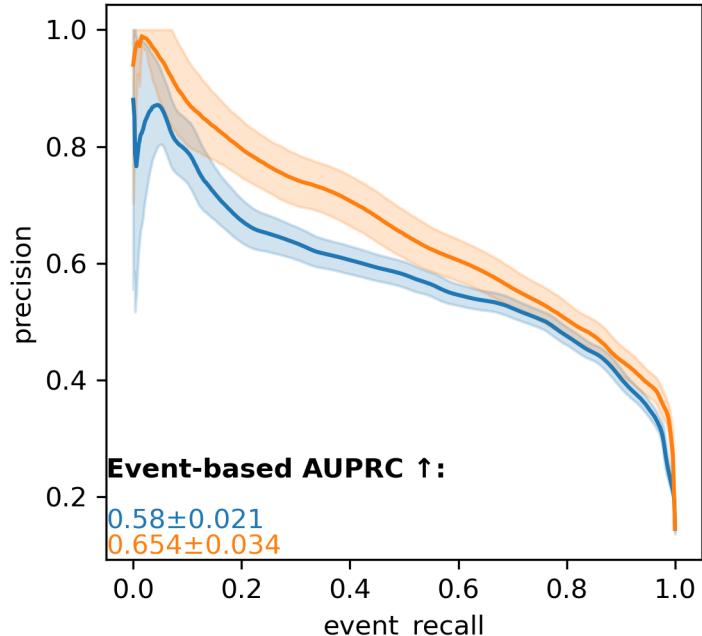
Precision / recall curve



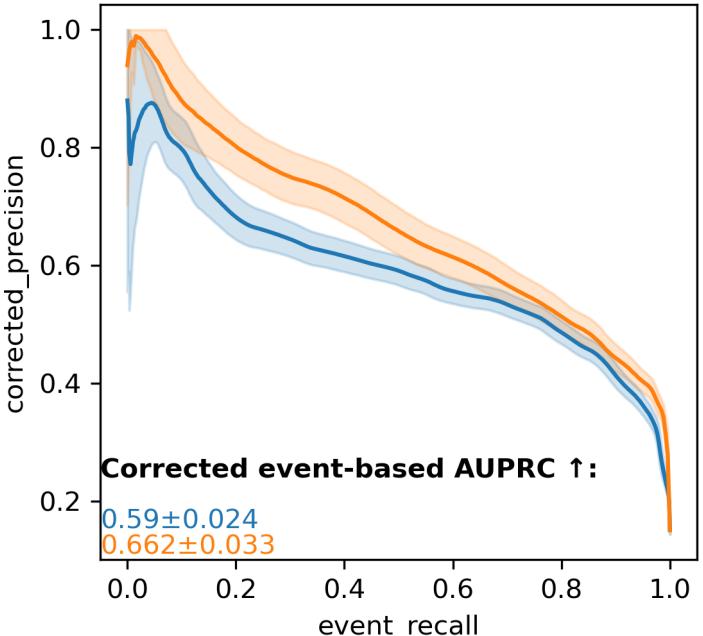
Corrected precision / recall curve



Precision / event-based recall curve



Corrected precision / event-based recall curve



Time Gap Analysis

Goal: Checking whether the time gap between the first correct alarm and the start of the corresponding event are similar across cohorts of patients

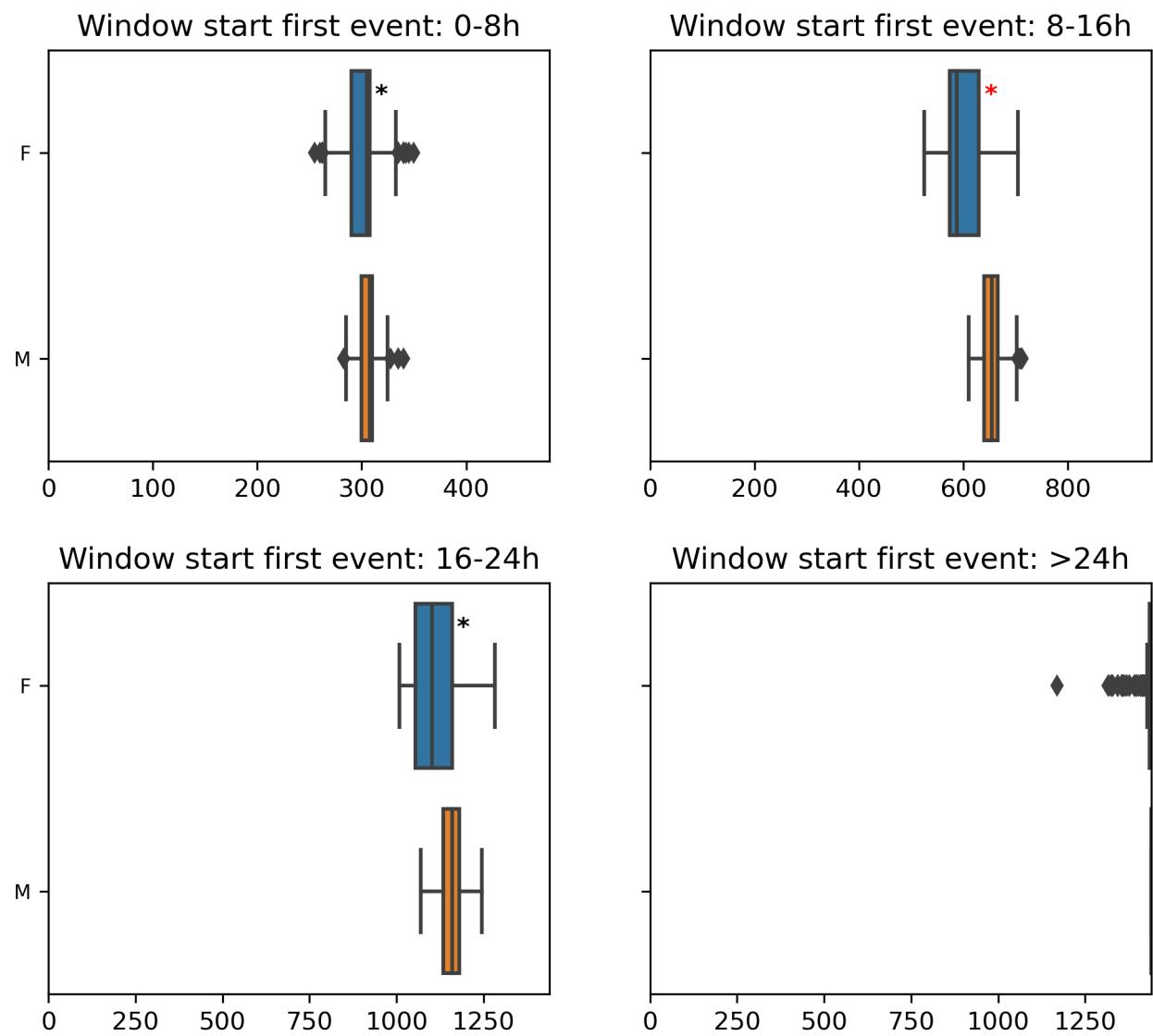
In the tables presenting the results of the statistical analysis, we present only window of start event and groups with a significant p-value (smaller than 0.001/nb_comparison) and whose delta is bigger than 0.

Top 3 categories with biggest deltas in terms of time gap discrepancy

Start event	Cat 1 (Δ in min)	Cat 2 (Δ in min)	Cat 3 (Δ in min)
0-8h	emergency (15.0)	Cardiovascular (10.0)	50-65 (8.75)
8-16h	F (67.5)	>85 (67.5)	Gastrointestinal (67.5)
16-24h	Metabolic (1132.5)	Gastrointestinal (100.0)	elective (85.0)
>24h	Metabolic (67.5)	>85 (42.5)	Trauma (30.0)

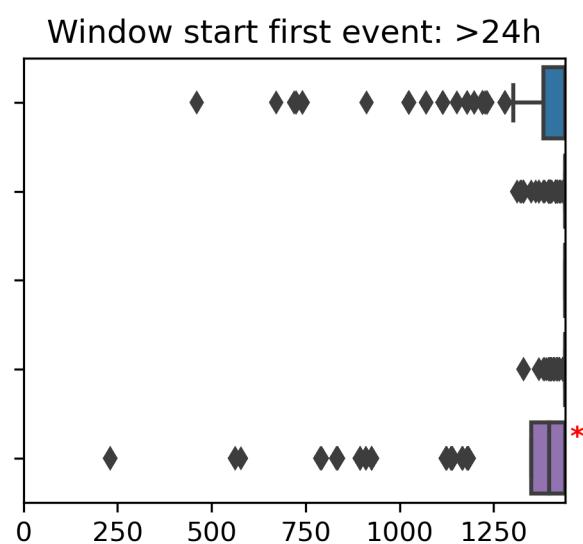
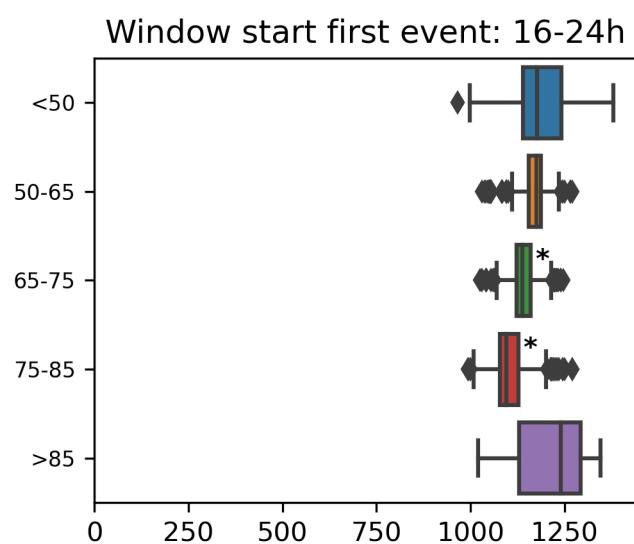
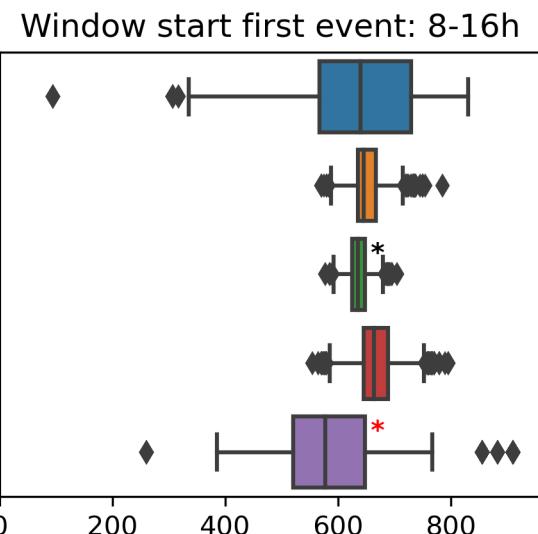
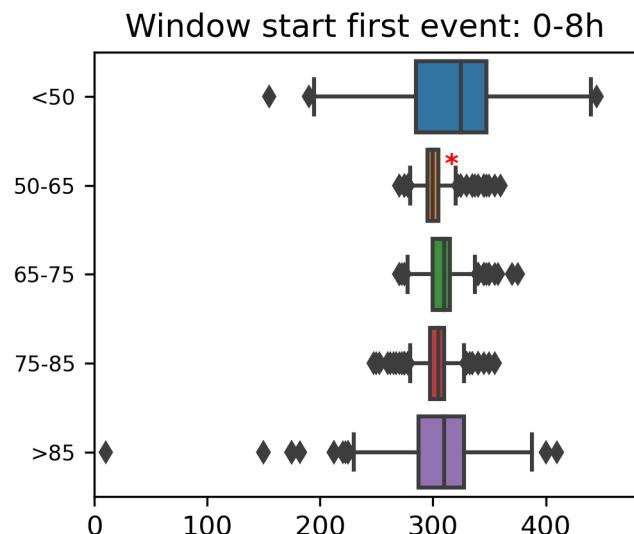
The following box plots show the median time gap between alarm and event for the different categories of patients depending on the time of stay when the event began.

Grouping by sex



Start event	Group with worst dist.	P-value	Delta (in min)
0-8h	F	1.55e-40	2.5
8-16h	F	1.06e-241	67.5
16-24h	F	1.29e-111	57.5

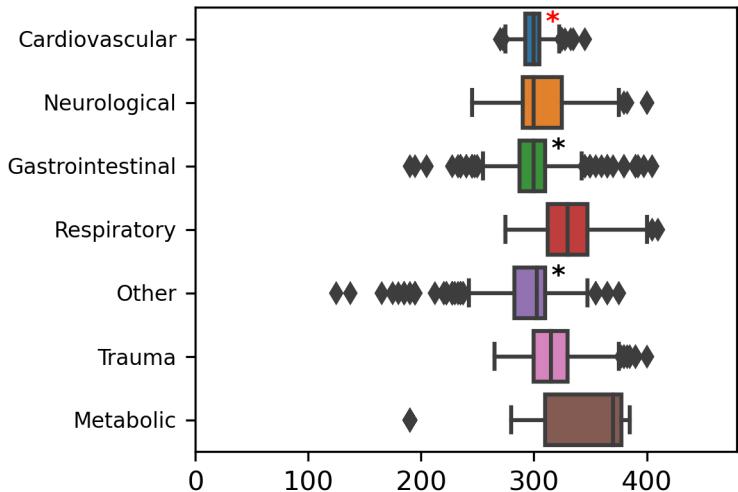
Grouping by age_group



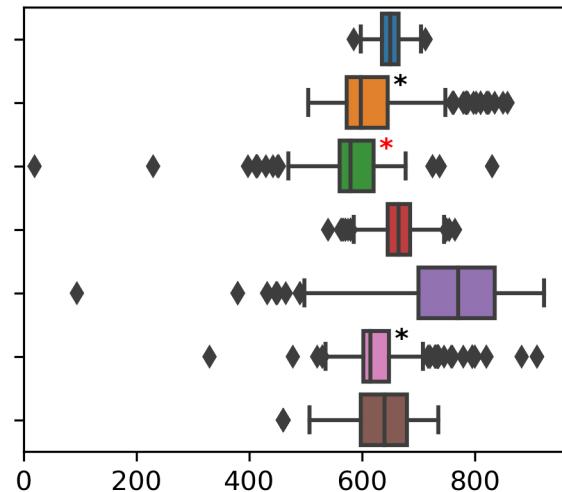
Start event	Group name	Group vs. rest	P-value	Delta (in min)
0-8h	<50	better	5.33e-15	20.0
0-8h	50-65	worse	1.78e-62	8.75
0-8h	65-75	better	1.40e-35	5.0
8-16h	50-65	better	3.95e-19	5.0
8-16h	65-75	worse	1.47e-51	10.0
8-16h	75-85	better	7.28e-98	26.25
8-16h	>85	worse	3.79e-79	67.5
16-24h	<50	better	6.26e-51	32.5
16-24h	50-65	better	3.24e-122	45.0
16-24h	65-75	worse	1.53e-22	32.5
16-24h	75-85	worse	7.40e-142	67.5
16-24h	>85	better	1.23e-59	92.5
>24h	>85	worse	1.87e-218	42.5

Grouping by APACHE_group

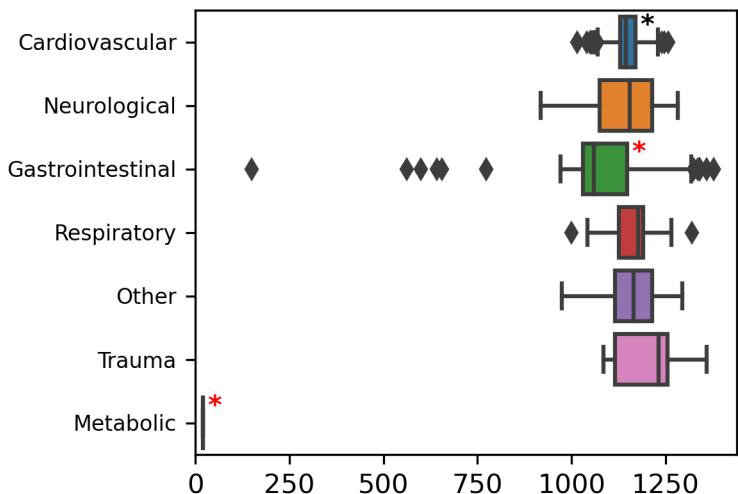
Window start first event: 0-8h



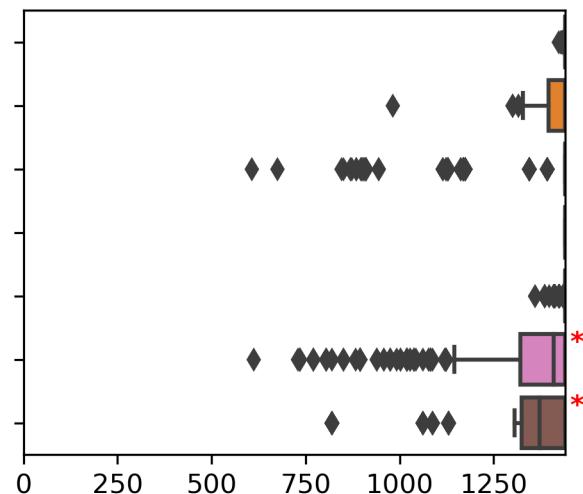
Window start first event: 8-16h



Window start first event: 16-24h



Window start first event: >24h

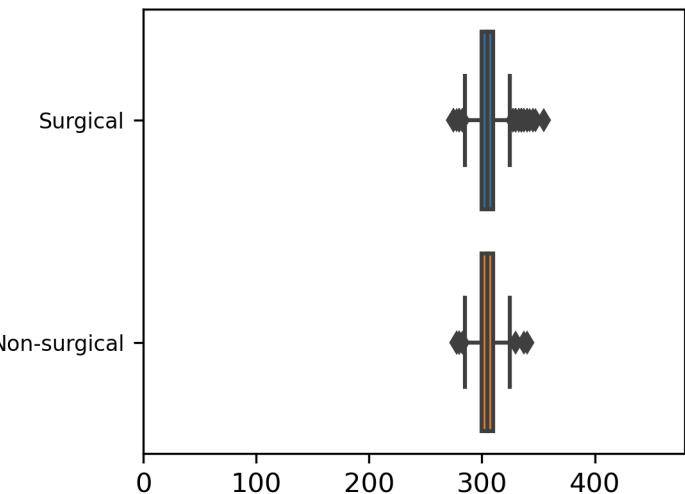


Start event	Group name	Group vs. rest	P-value	Delta (in min)
0-8h	Cardiovascular	worse	3.02e-141	10.0
0-8h	Gastrointestinal	worse	8.52e-37	5.0
0-8h	Respiratory	better	5.85e-223	30.0
0-8h	Other	worse	3.59e-22	2.5
0-8h	Trauma	better	1.22e-44	10.0
0-8h	Metabolic	better	8.07e-73	65.0
8-16h	Cardiovascular	better	1.59e-46	12.5
8-16h	Neurological	worse	2.48e-101	47.5
8-16h	Gastrointestinal	worse	3.40e-275	67.5
8-16h	Respiratory	better	3.14e-116	25.0
8-16h	Other	better	6.41e-160	130.0
8-16h	Trauma	worse	2.84e-72	30.0
16-24h	Cardiovascular	worse	1.7e-05	10.0
16-24h	Gastrointestinal	worse	2.65e-107	100.0
16-24h	Respiratory	better	3.97e-50	33.75
16-24h	Other	better	1.04e-18	17.5
16-24h	Trauma	better	9.57e-61	81.25

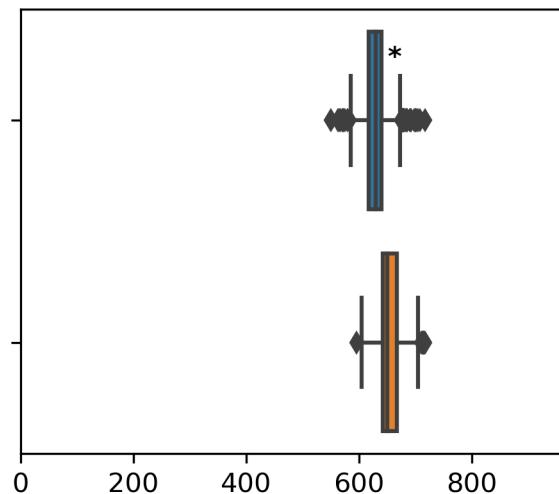
16-24h	Metabolic	worse	5.42e-184	1132.5
>24h	Trauma	worse	3.20e-215	30.0
>24h	Metabolic	worse	1.18e-201	67.5

Grouping by surgical_status

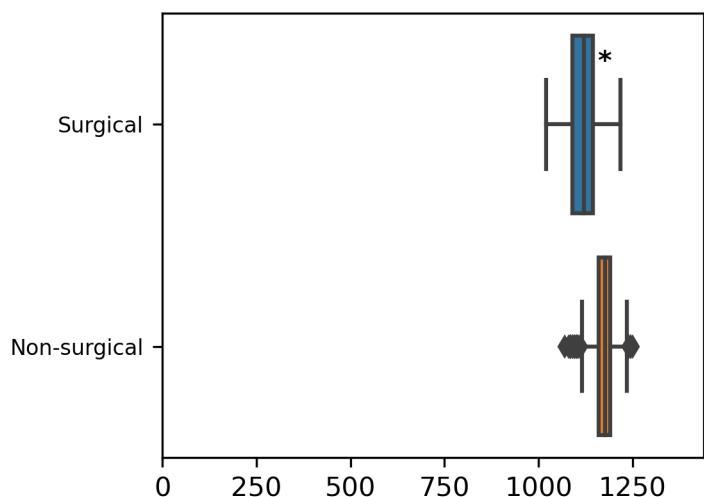
Window start first event: 0-8h



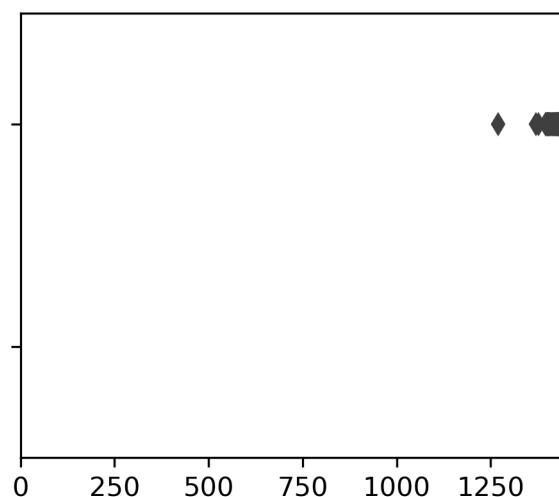
Window start first event: 8-16h



Window start first event: 16-24h



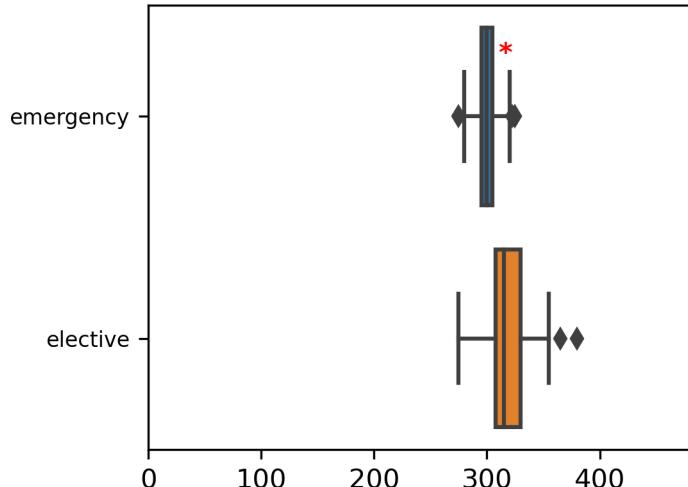
Window start first event: >24h



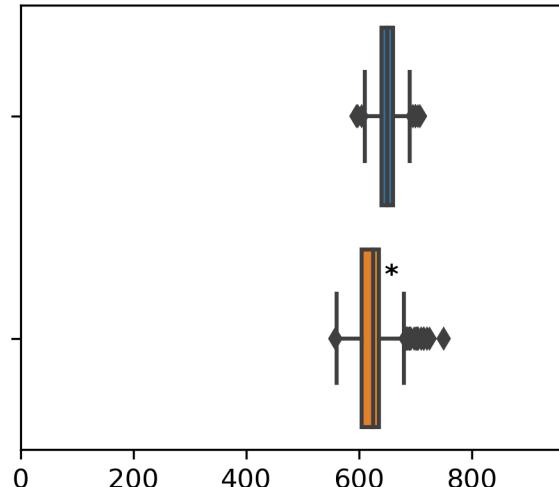
Start event	Group with worst dist.	P-value	Delta (in min)
8-16h	Surgical	2.33e-132	20.0
16-24h	Surgical	1.37e-185	57.5

Grouping by admission_type

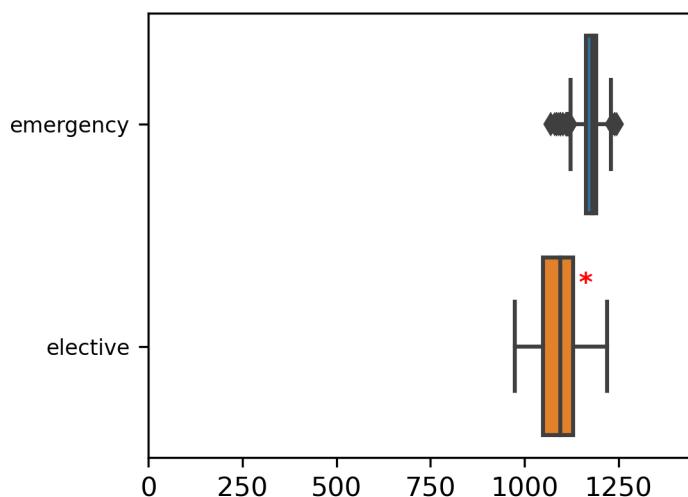
Window start first event: 0-8h



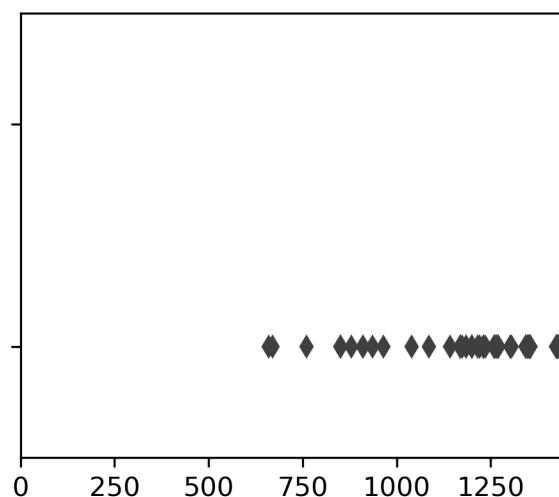
Window start first event: 8-16h



Window start first event: 16-24h



Window start first event: >24h



Start event	Group with worst dist.	P-value	Delta (in min)
0-8h	emergency	1.29e-189	15.0
8-16h	elective	1.43e-150	25.0
16-24h	elective	1.60e-241	85.0