# Class-Dependent Miscalibration Severely Degrades Selective Prediction in Multimodal Clinical Prediction Models

**L. Julián Lechuga López[1,2]**
**Farah E. Shamout[1,2]**
**Tim G. J. Rudner[3]**

LEOPOLDO.LECHUGA@NYU.EDU
FARAH.SHAMOUT@NYU.EDU
TIM.RUDNER@UTORONTO.CA

[1] *NYU Abu Dhabi, Abu Dhabi, UAE*

[2] *New York University, New York, NY, USA*

[3] *University of Toronto, Toronto, ON, Canada*

## Abstract

As artificial intelligence systems transition from research to clinical deployment, ensuring their reliability becomes critical for clinical decision-making tasks, as incorrect predictions can have serious consequences. Deploying AI in healthcare therefore requires prediction systems with robust safeguards against error, such as selective prediction, where uncertain predictions are deferred to human experts for review. In this study, we carefully evaluate the reliability of uncertainty-based selective prediction for multi-label clinical condition classification using multimodal data. Our findings show that models often exhibit severe class-dependent miscalibration causing predictive performance to *degrade* under uncertainty-guided selective prediction—attributing high uncertainty to correct predictions and low uncertainty to incorrect predictions. Our evaluation highlights fundamental shortcomings of commonly used evaluation metrics for clinical AI. To address these shortcomings, we propose practical recommendations for calibration-aware model assessment and selective prediction design, offering a pathway to safer, more reliable AI systems that clinicians and patients can trust.

**Keywords:** Calibration; Multimodal learning; Selective prediction.

**Data and Code Availability** This study uses publicly available benchmark datasets: MIMIC-IV Johnson et al. (2021) and MIMIC-CXR Johnson et al. (2019). Both datasets are available to researchers through the PhysioNet platform, subject to credentialed access and data use agreements. We plan to make our code available in a future release.

**Institutional Review Board (IRB)** This research uses only publicly available, de-identified datasets (MIMIC-IV and MIMIC-CXR) that do not constitute human subjects research. Therefore, IRB approval was not required.

## 1. Introduction

Machine learning is increasingly embedded into the healthcare sector to support clinical decision-making, improve diagnostic performance, accelerate drug discovery, and optimize patient management (Paul et al., 2025; Hanna et al., 2025). From enhancing disease detection (Mall et al., 2022) to personalizing treatment plans (Agarwal, 2024), machine learning-driven systems have demonstrated transformative potential. However, in high-stakes clinical environments, exhibiting good performance, such as high accuracy, is not sufficient. Models in healthcare must provide fail-safe mechanisms to ensure they are reliable and interpretable, to allow clinicians to recognize when predictions are unreliable so they can intervene accordingly (Javed et al., 2024).

One such mechanism is selective prediction, where a model can abstain from making a prediction and request a review by a human expert if its predictive uncertainty is high (Geifman and El-Yaniv, 2017). By abstaining in cases of high uncertainty, selective prediction provides a fail-safe mechanism against critical errors that can endanger patients. However, for selective prediction to improve model safety and robustness, models must be able to provide reliable uncertainty estimates. If uncertainty estimates are systematically miscalibrated, selective prediction may fail to identify risky cases, undermining its value in clinical decision-making.

Unfortunately, evidence from unimodal domains such as computer vision and time-series analysis shows that many models systematically misestimate uncertainty, often being overconfident when wrong (Guo et al., 2017; Morey et al., 2025; Deng et al., 2025). In healthcare settings, these challenges may be further compounded when having to deal with multimodal data, where heterogeneous sources such as electronic health records (EHR), imaging, and text are integrated to provide a holistic view of a patient's state (Simon et al., 2025). While multimodal methods frequently improve discrimination metrics, their calibration properties, and thus their reliability for selective prediction, remain poorly understood (Zhao et al., 2024).

In this paper, we study whether state-of-the-art multimodal clinical condition classification models provide reliable uncertainty estimates that enable effective selective prediction. We find that although multimodal models outperform unimodal baselines on standard evaluation metrics, multimodal fusion substantially degrades selective prediction, calling into question its suitability as a fail-safe mechanism in safety-critical healthcare settings (Table 1).

To understand the roots of this failure, we perform a careful analysis across unimodal and multimodal models trained on one of the largest publicly-available multimodal datasets, consisting of EHR and chest X-ray data. This comparison reveals that multimodal fusion does not consistently improve calibration and, in many cases, exacerbates class-dependent miscalibration. In particular, we consider class-dependent expected calibration errors and find that positive-class predictions are systematically overconfident, especially for conditions underrepresented in the data, undermining the effectiveness of selective prediction.

Finally, we investigate whether simple loss-upweight training strategies can mitigate these issues and find that they offer only modest improvements, falling short of resolving class-dependent miscalibration in multimodal models. Our findings indicate that state-of-the-art multimodal models cannot deliver selective prediction as a reliable fail-safe mechanism, underscoring calibration-aware selective prediction as a key open challenge for safe deployment.

To summarize, our key contributions are:

1. We show that selective prediction in multimodal clinical condition classification models degrades performance across a wide range of clinical conditions.

2. We investigate the causes of this failure by comparing unimodal and multimodal models and show that the degradation in selective prediction is driven by class-dependent miscalibration.

3. We introduce a simple label-dependent loss up-weighting scheme that partially corrects class-dependent miscalibration and improves selective prediction.

4. We quantify correlation changes in selective AUC with stratified ECE AUC revealing that, substantial miscalibration remains across architectures, underscoring the need for more targeted approaches for safe multimodal clinical condition classification.

We find that high predictive performance in multimodal clinical condition classification often coincides with severe miscalibration, especially for rare conditions. In safety-critical settings, a model that looks accurate but fails to flag its own errors cannot be trusted. Although selective prediction could enable safety guardrails, current multimodal models systematically undermine this promise. Closing this gap demands models that couple strong discrimination with reliable selective prediction, making calibration-aware design essential for trustworthy clinical AI.

## 2. Related Work

Multimodal learning is increasingly used in clinical machine learning for integrating heterogeneous data sources such as structured EHR data, medical imaging, and free-text clinical reports (Warner et al., 2024; Simon et al., 2025). By combining complementary modalities, these models aim to emulate how clinicians synthesize diverse information for diagnosis and treatment planning (Ng et al., 2024; Lin et al., 2025).

Multimodal learning has been used in varying contexts: in some cases referring to multiple imaging views of the same modality (e.g., frontal and lateral chest X-rays) (Warner et al., 2024), in others describing fusion across fundamentally different data types such as imaging and tabular or textual inputs (Niu et al., 2024; Lee et al., 2024). While earlier work often focused on combining multiple imaging modalities for classification or segmentation (Sangeetha et al., 2024; Adebiyi et al., 2024), recent advances in large language models and vision–language models have broadened multimodal clinical systems to include combinations of radiology images, structured
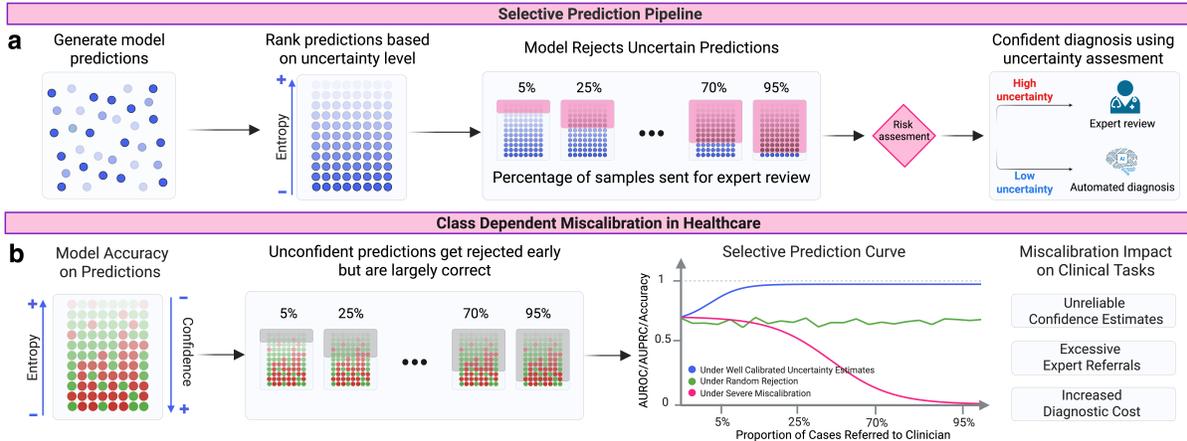
Figure 1: **Selective Prediction Can Serve as a Fail-Safe Mechanism in Safety-Critical Healthcare Settings.** **a)** Selective prediction pipeline: a model outputs prediction confidence/uncertainty scores; low-confidence cases are rejected and referred to an expert, yielding coverage-controlled evaluation and diagnosis. **b)** Calibration drives selective behavior: miscalibration causes predictions that are mostly correct to have high uncertainty (low confidence) and incorrect predictions to have high confidence. Theoretical selective prediction curves show the difference between a well-calibrated model, random selection, and a severely miscalibrated model, showing improving, flat, and degrading performance with increasing coverage. Clinically, overconfidence risks missed or delayed diagnoses, while underconfidence triggers unnecessary rejections and shifts excess workload to clinicians.

vitals, free-text notes, and even patient–provider conversations (Mahesh et al., 2024; Zambrano Chaves et al., 2025). This highlights the growing importance of multimodal integration in clinical AI, for tasks such as clinical question answering (Shoham and Rappoport, 2024), automated report generation (Chen et al., 2024), and medical documentation (Khan et al., 2025), underscoring the need to assess not only discrimination performance but also calibration and reliability when multiple heterogeneous sources are fused. Despite this progress, prior work in multimodal clinical ML has paid little attention to calibration, leaving open questions about whether improvements in predictive performance translate into reliable probability estimates for safe deployment (Zhao et al., 2024; Deng et al., 2025).

To investigate these challenges in a concrete and clinically meaningful setting, we focus on multimodal fusion of EHR time-series and chest X-ray (CXR) imaging, which is particularly relevant for critical care tasks such as diagnosis, mortality prediction, and intervention planning (Schilcher et al., 2024). Fusion architectures, such as MedFuse (Hayat et al., 2022), DrFuse (Yao et al., 2024) or Metra (Khader et al.,

2023), have shown that even relatively simple strategies like early or late concatenation can yield performance gains on benchmarks such as MIMIC (Johnson et al., 2021). This paired EHR–CXR setting provides a practical and widely used benchmark for investigating how multimodal learning impacts selective prediction, calibration, and reliability in a clinical setting (Figure 1.b). In Appendix 3 we provide the theoretical details of our selective prediction pipeline and uncertainty quantification procedure.

## 3. Methods

**Clinical Condition Classification.** The objective of this multilabel classification task is to predict whether a patient presents any of 25 chronic, acute, or mixed clinical conditions during a given stay in the intensive care unit (ICU) (Hayat et al., 2022). Derived from one of the largest publicly-available multimodal datasets, this task has become a standard benchmark for multimodal predictive modeling in healthcare. The clinical conditions and their prevalence are detailed in Appendix 1.
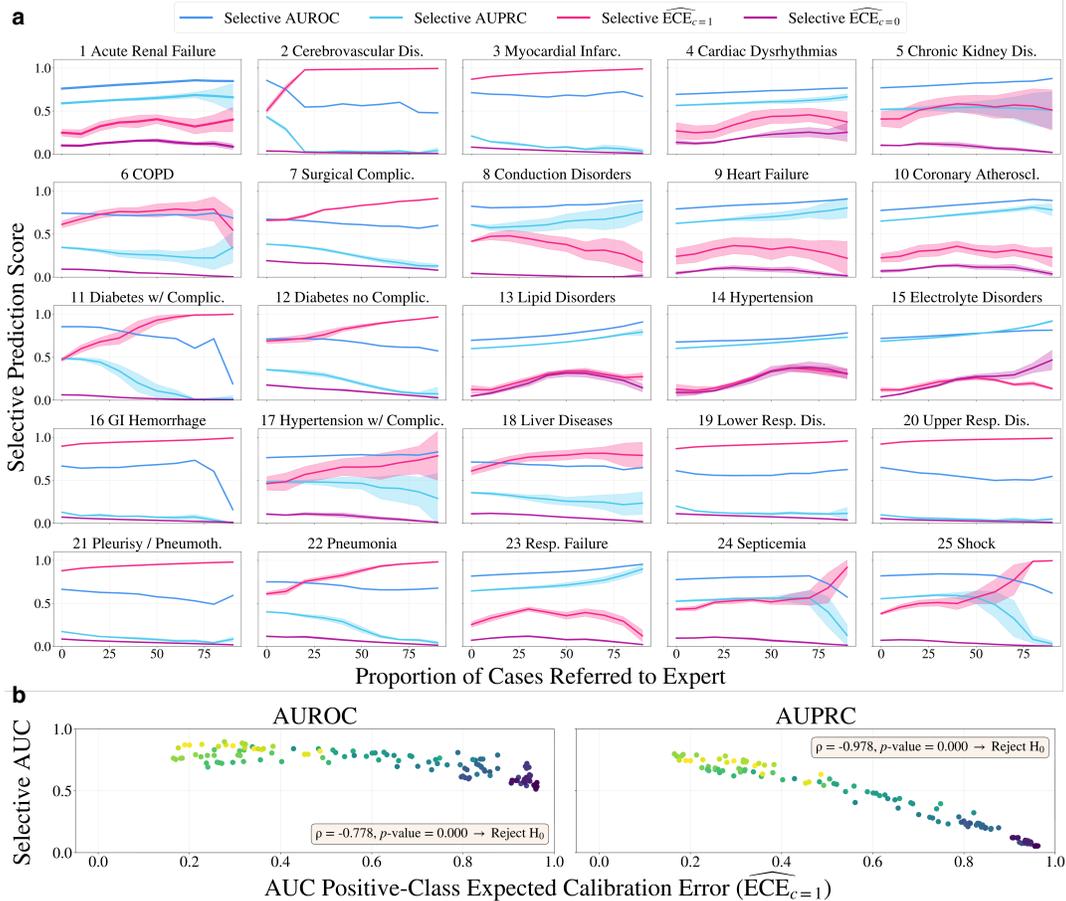
Figure 2: **How Does Calibration Shape Selective Prediction? a)** Selective AUROC, AUPRC, and class-stratified ECE across 25 conditions showing how performance evolves as low-confidence cases are rejected highlighting the imbalance between calibration and discrimination. Shaded regions show variance over five seeds. **b)** Higher positive-class ECE AUC consistently predicts lower selective AUC across all conditions, showing that overconfidence in positive predictions drives instability in selective evaluation. The null hypothesis of no non-negative correlation between ECE and Selective AUC across conditions, is rejected for both metrics (Spearman's rank $H_0$, $p < 0.05$).

**Multimodal Dataset.** We construct paired multimodal samples from MIMIC-IV (Johnson et al., 2021) (structured EHR time-series) and MIMIC-CXR (Johnson et al., 2019) (frontal-view chest X-rays), such that each patient sample contains both modalities: $x^{ehr}$, $x^{cxr}$ and multilabel ground truth $y_{ehr}$. We follow standard preprocessing steps and do not allow overlap of patients across training and test splits, as in previous work (Hayat et al., 2022).

### 3.1. Model Variants

**Architecture.** We adopt the MedFuse architecture (Hayat et al., 2022) to train a multimodal

classifier $f(\cdot)$ on paired EHR and CXR data $\mathcal{D} = \{(x_n^{ehr}, x_n^{cxr}, y_n)\}_{n=1}^{N}$. EHR time-series are encoded using a two-layer LSTM ($\Phi_{ehr}$) (Hochreiter and Schmidhuber, 1997), and CXR images using a ResNet-34 ($\Phi_{cxr}$) (He et al., 2015), with latent representations concatenated and passed to a classification head $g(\cdot)$ with sigmoid outputs. Models are trained with the standard binary cross-entropy loss, i.e. $\log p(y|\hat{y}) = -(y \log \hat{y} + (1 - y) \log(1 - \hat{y}))$.

**Unimodal Baselines.** We train unimodal models using only $\Phi_{ehr}$ (unimodal EHR) and $\Phi_{cxr}$ (unimodal CXR) to establish reference points for performance and calibration. The CXR encoder is initial-

4

Table 1: **Performance Across Conditions.** Aggregate metrics suggest that multimodal architectures outperform unimodal baselines. However, they mask substantial condition-specific variability, particularly in calibration. (**Best**, <u>Second Best</u>).

| Model | AUROC | AUPRC | Selective AUROC | Selective AUPRC | $\overline{\text{ECE}} \downarrow$ |
|---|---|---|---|---|---|
| EHR (LSTM) (Hochreiter and Schmidhuber, 1997) | $0.681_{\pm0.007}$ | $0.371_{\pm0.010}$ | $0.625_{\pm0.115}$ | $0.292_{\pm0.091}$ | $\mathbf{1.42_{\pm0.325}}$ |
| CXR (ResNet-34) (He et al., 2015) | $0.654_{\pm0.013}$ | $0.350_{\pm0.012}$ | $0.643_{\pm0.085}$ | $0.342_{\pm0.077}$ | $9.41_{\pm5.826}$ |
| MedFuse (Hayat et al., 2022) | $0.739_{\pm0.005}$ | $0.449_{\pm0.008}$ | $0.705_{\pm0.136}$ | $0.396_{\pm0.112}$ | $2.08_{\pm0.725}$ |
| MedFuse (Loss Upweight) (Hayat et al., 2022) | $\mathbf{0.742_{\pm0.006}}$ | $\mathbf{0.447_{\pm0.010}}$ | $\mathbf{0.764_{\pm0.098}}$ | $\mathbf{0.532_{\pm0.002}}$ | $4.88_{\pm1.470}$ |
| DrFuse (Yao et al., 2024) | $0.726_{\pm0.007}$ | $0.418_{\pm0.010}$ | $0.677_{\pm0.139}$ | $0.352_{\pm0.116}$ | $2.54_{\pm1.15}$ |
| MeTra (Khader et al., 2023) | $0.707_{\pm0.019}$ | $0.399_{\pm0.015}$ | $0.661_{\pm0.125}$ | $0.333_{\pm0.101}$ | $2.32_{\pm0.95}$ |
| Δ MedFuse vs EHR | +0.058 (8.5%) | +0.078 (21.0%) | +0.080 (12.8%) | +0.104 (35.6%) | +0.66 (46%) |
| Δ MedFuse (Loss Upweight) vs EHR | +0.061 (9.0%) | +0.076 (20.5%) | +0.139 (22.2%) | +0.240 (82.2%) | +3.46 (243.7%) |

ized with ImageNet weights, and the EHR encoder is trained from scratch.

**Multimodal Architectures.** The original Med-Fuse fusion model serves as our deterministic multimodal baseline, allowing us to measure the calibration impact of modality fusion and providing a reference for uncertainty-aware variants. To assess the robustness of our findings to the choice of architecture, we additionally evaluate two alternative multimodal backbones: DrFuse (Yao et al., 2024), which learns EHR and CXR representations with divergence-based alignment, and MeTra (Khader et al., 2023), which encodes images and clinical variables using a transformer-based cross-modal fusion encoder.

## 4. How Does Calibration Shape Selective Prediction?

**Calibration Drives Selective Performance.** We next examine selective prediction thresholds for all 25 conditions using MedFuse, plotting selective AUROC and AUPRC alongside stratified ECE (Figure 2.a). Ideally, selective metrics should increase as uncertain predictions are correctly rejected, but we observe that many conditions deviate from this expectation. At extreme thresholds, selective AUROC and AUPRC collapse when the class distribution becomes particularly imbalanced. The mid-threshold regimes reveal more concerning behaviors where selective metrics stagnate or degrade, indicating that the model does not consistently reject the "right" cases.

Linking these failures to calibration reveals clear patterns. Conditions with low, balanced positive and negative ECE exhibit stable and improving selective AUROC/AUPRC (e.g., Labels 13 to 15). By contrast, high positive-class ECE produces noisy, unstable selective curves: AUROC hovers near chance and

AUPRC declines sharply (e.g., Labels 2, 22, 24 and 25). Variance across random seeds follows the same trend, with poor calibration associated with wide, unstable spreads.

**Calibration as a Predictor of Evaluation Stability.** Across conditions (Figure 2.b), positive-class ECE AUC and Selective AUC are consistently negatively correlated: higher positive-class ECE predicts smaller or even negative changes in selective prediction metrics. The effect is strongest for AUPRC, where degradation is nearly monotonic, while AUROC shows greater variability but follows the same trend. We fomally quantify this with the null hypothesis of no non-negative association between ECE and Selective AUC across conditions, which is rejected for both metrics (Spearman's rank correlation, $p < 0.05$).

Overall, these findings establish positive-class calibration as a reliable predictor of selective prediction stability. High positive-class ECE serves as an early warning that selective prediction will fail to yield meaningful improvements over standard evaluation. Since this holds across all conditions, it reflects a generalizable phenomenon rather than a label-specific artifact.

> **Takeaway.** Positive-class calibration is a leading indicator of selective reliability: when $\text{ECE}_{c=1}$ is high, selective AUROC/AUPRC fail to improve and can deteriorate.

**Extended Results.** In Appendix 4 we complement this analysis with full quantitative comparisons on individual modalities (CXR, EHR) and multimodal architectures (MedFuse, DrFuse and MeTra) across different evaluation settings.

## 5. Conclusion

Our empirical findings highlight a central challenge: *multimodal fusion improves discrimination but undermines calibration, making selective prediction unreliable.* Addressing this gap requires moving beyond aggregate performance and developing calibration-aware selective prediction methods that can provide trustworthy fail-safe mechanisms. Without this, the deployment of multimodal clinical AI risks offering the illusion of safety while leaving critical vulnerabilities unaddressed.

# References

Abdulmateen Adebiyi, Nader Abdalnabi, Emily Hoffman Smith, Jesse Hirner, Eduardo J Simoes, Mirna Becevic, and Praveen Rao. Accurate skin lesion classification using multimodal learning on the ham10000 dataset. *MedRxiv*, pages 2024–05, 2024.

Shashank Agarwal. Machine learning based personalized treatment plans for chronic conditions. In *2024 2nd International Conference on Intelligent Data Communication Technologies and Internet of Things (IDCIoT)*, pages 1127–1132. IEEE, 2024.

Zhihong Chen, Maya Varma, Jean-Benoit Delbrouck, Magdalini Paschali, Louis Blankemeier, Dave Van Veen, Jeya Maria Jose Valanarasu, Alaa Youssef, Joseph Paul Cohen, Eduardo Pontes Reis, et al. Chexagent: Towards a foundation model for chest x-ray interpretation. *arXiv preprint arXiv:2401.12208*, 2024.

Jiawen Deng, Mohamed E Elghobashy, Kathleen Zang, Shubh K Patel, Eddie Guo, and Kiyan Heybati. So you've got a high auc, now what? an overview of important considerations when bringing machine-learning models from computer to bedside. *Medical Decision Making*, page 0272989X251343082, 2025.

Stefan Depeweg, Jose-Miguel Hernandez-Lobato, Finale Doshi-Velez, and Steffen Udluft. Decomposition of uncertainty in bayesian deep learning for efficient and risk-sensitive learning. In *International conference on machine learning*, pages 1184–1193. PMLR, 2018.

Ran El-Yaniv et al. On the foundations of noise-free selective classification. *Journal of Machine Learning Research*, 11(5), 2010.

Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. *Advances in neural information processing systems*, 30, 2017.

Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pages 1321–1330. PMLR, 2017.

Matthew G Hanna, Liron Pantanowitz, Rajesh Dash, James H Harrison, Mustafa Deebajah, Joshua Pantanowitz, and Hooman H Rashidi. Future of artificial intelligence (ai)-machine learning (ml) trends in pathology and medicine. *Modern Pathology*, page 100705, 2025.

Nasir Hayat, Krzysztof J. Geras, and Farah E. Shamout. Medfuse: Multi-modal fusion with clinical time-series data and chest x-ray images. In *Proceedings of the 7th Machine Learning for Healthcare Conference*, volume 182 of *Proceedings of Machine Learning Research*, pages 479–503. PMLR, August 2022.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.

Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.

Haseeb Javed, Shaker El-Sappagh, and Tamer Abuhmed. Robustness in deep learning models for medical diagnostics: security and adversarial challenges towards robust ai applications. *Artificial Intelligence Review*, 58(1):12, 2024.

Alistair Johnson, Lucas Bulgarelli, Tom Pollard, Leo Anthony Celi, Roger Mark, and S Horng IV. Mimic-iv-ed. *PhysioNet*, 2021.

Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Roger G Mark, and Steven Horng. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1):317, 2019.

Firas Khader, Jakob Nikolas Kather, Gustav Müller-Franzes, Tianci Wang, Tianyu Han, Soroosh Tayebi Arasteh, Karim Hamesch, Keno Bressem, Christoph Haarburger, Johannes Stegmaier, et al. Medical transformer for multimodal survival prediction in intensive care: integration of imaging and non-imaging data. *Scientific Reports*, 13(1): 10666, 2023.

Wasif Khan, Seowung Leem, Kyle B See, Joshua K Wong, Shaoting Zhang, and Ruogu Fang. A comprehensive survey of foundation models in medicine. *IEEE Reviews in Biomedical Engineering*, 2025.

Chih-Kuo Lee, Ting-Li Chen, Jeng-En Wu, Min-Tsun Liao, Chiehhung Wang, Weichung Wang, and Cheng-Ying Chou. Multimodal deep learning models utilizing chest x-ray and electronic health record

data for predictive screening of acute heart failure in emergency department. *Computer Methods and Programs in Biomedicine*, 255:108357, 2024.

Qika Lin, Yifan Zhu, Xin Mei, Ling Huang, Jingying Ma, Kai He, Zhen Peng, Erik Cambria, and Mengling Feng. Has multimodal learning delivered universal intelligence in healthcare? a comprehensive survey. *Information Fusion*, 116:102795, 2025.

Nandhini Mahesh, Chitralekha S Devishamani, Keerthana Raghu, Maanasi Mahalingam, Pragathi Bysani, Arjun V Chakravarthy, and Rajiv Raman. Advancing healthcare: the role and impact of ai and foundation models. *American Journal of Translational Research*, 16(6):2166, 2024.

Shachi Mall, Ashutosh Srivastava, Bireshwar Dass Mazumdar, Manmohan Mishra, Sunil L Bangare, and A Deepak. Implementation of machine learning techniques for disease diagnosis. *Materials Today: Proceedings*, 51:2198–2201, 2022.

Dane A Morey, Michael F Rayo, and David D Woods. Empirically derived evaluation requirements for responsible deployments of ai in safety-critical settings. *npj Digital Medicine*, 8(1):374, 2025.

Jeremy Y Ng, Holger Cramer, Myeong Soo Lee, and David Moher. Traditional, complementary, and integrative medicine and artificial intelligence: Novel opportunities in healthcare. *Integrative medicine research*, 13(1):101024, 2024.

Shuai Niu, Jing Ma, Liang Bai, Zhihua Wang, Li Guo, and Xian Yang. Ehr-knowgen: Knowledge-enhanced multimodal learning for disease diagnosis generation. *Information Fusion*, 102:102069, 2024.

Jibon Kumar Paul, Mahir Azmal, Omar Faruk Talukder, ANM Shah Newaz Been Haque, Meghla Meem, and Ajit Ghosh. Harnessing machine learning for improved diagnosis, drug discovery, and patient care. *Computational and Structural Biotechnology Reports*, page 100051, 2025.

Stephan Rabanser and Nicolas Papernot. What does it take to build a performant selective classifier? *arXiv preprint arXiv:2510.20242*, 2025.

Stephan Rabanser, Anvith Thudi, Kimia Hamidieh, Adam Dziedzic, and Nicolas Papernot. Selective classification via neural network training dynamics. *arXiv preprint arXiv:2205.13532*, 2022.

Tim G. J. Rudner, Sanyam Kapoor, Shikai Qiu, and Andrew Gordon Wilson. Function-space regularization in neural networks: A probabilistic perspective. In *Proceedings of the 40th International Conference on Machine Learning*, Proceedings of Machine Learning Research. PMLR, 2023a.

Tim G. J. Rudner, Sanyam Kapoor, Shikai Qiu, and Andrew Gordon Wilson. Should we learn most likely functions or parameters? In *Advances in Neural Information Processing Systems 36*, 2023b.

SKB Sangeetha, Sandeep Kumar Mathivanan, P Karthikeyan, Hariharan Rajadurai, Basu Dev Shivahare, Saurav Mallik, and Hong Qin. An enhanced multimodal fusion deep learning neural network for lung cancer classification. *Systems and Soft Computing*, 6:200068, 2024.

Jörg Schilcher, Alva Nilsson, Oliver Andlid, and Anders Eklund. Fusion of electronic health records and radiographic images for a multimodal deep learning prediction model of atypical femur fractures. *Computers in Biology and Medicine*, 168:107704, 2024.

Ofir Ben Shoham and Nadav Rappoport. Medconceptsqa: Open source medical concepts qa benchmark. *Computers in Biology and Medicine*, 182:109089, 2024.

Benjamin D Simon, Kutsev Bengisu Ozyoruk, David G Gelikman, Stephanie A Harmon, and Barış Türkbey. The future of multimodal artificial intelligence models for integrating imaging and clinical metadata: a narrative review. *Diagnostic and Interventional Radiology*, 31(4):303, 2025.

Elisa Warner, Joonsang Lee, William Hsu, Tanveer Syeda-Mahmood, Charles E Kahn Jr, Olivier Gevaert, and Arvind Rao. Multimodal machine learning in image-based and clinical biomedicine: Survey and prospects. *International journal of computer vision*, 132(9):3753–3769, 2024.

Wenfang Yao, Kejing Yin, William K Cheung, Jia Liu, and Jing Qin. Drfuse: Learning disentangled representation for clinical multi-modal fusion with missing modality and modal inconsistency. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pages 16416–16424, 2024.

Juan Manuel Zambrano Chaves, Shih-Cheng Huang, Yanbo Xu, Hanwen Xu, Naoto Usuyama, Sheng

Zhang, Fei Wang, Yujia Xie, Mahmoud Khademi, Ziyi Yang, et al. A clinically accessible small multimodal radiology model and evaluation metric for chest x-ray findings. *Nature Communications*, 16 (1):3108, 2025.

Tianyi Zhao, Liangliang Zhang, Yao Ma, and Lu Cheng. A survey on safe multi-modal learning systems. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 6655–6665, 2024.

# Appendix

**Table of Contents**

## Reproducibility

We have made a significant effort to ensure the reproducibility of our results. An anonymized implementation of our method is provided at anonymous.4open.science/r/medcalibration-B187, which includes training, evaluation, and analysis scripts. The experimental setup—including hyperparameters, model configurations, and sampling parameters are described in Appendix 2. All datasets used in our experiments are publicly available, and we additionally provide scripts for data preparation. We provide model checkpoints for all models used in this paper with instructions for reproducing results in the `README.md` file in the released code.

# 1. Clinical Condition Dataset Distribution

Table 2: **Clinical Conditions and Dataset Characteristics.** Overview of the 25 conditions analyzed in this study, including their associated organ/system, primary diagnostic modality (CXR or EHR), and the prevalence of positive and negative samples in the training set.

| Clinical Condition | Organ | Modality CXR | Modality EHR | Train Set Prevalence Positive | Train Set Prevalence Negative |
|---|---|---|---|---|---|
| 1 Acute and unspecified renal failure | Kidneys | - | ✓ | 32.32% (2,532) | 67.68% (5,309) |
| 2 Acute cerebrovascular disease | Brain (Cerebrovascular system) | - | ✓ | 7.99% (625) | 92.01% (7,216) |
| 3 Acute myocardial infarction | Heart (Coronary arteries) | - | ✓ | 8.57% (675) | 91.43% (7,166) |
| 4 Cardiac dysrhythmias | Heart (Conduction system) | - | ✓ | 37.36% (2,942) | 62.64% (4,899) |
| 5 Chronic kidney disease | Kidneys | ✓ | ✓ | 23.34% (1,824) | 76.66% (6,017) |
| 6 Chronic obstructive pulmonary disease | Lungs (Airways) | ✓ | ✓ | 16.40% (1,269) | 83.60% (6,572) |
| 7 Complications of surgical procedures | Multiple systems (context-specific) | - | ✓ | 21.75% (1,703) | 78.25% (6,138) |
| 8 Conduction disorders | Heart (Electrical system) | ✓ | ✓ | 10.79% (838) | 89.21% (7,003) |
| 9 Congestive heart failure; nonhypertensive | Heart | - | ✓ | 31.03% (2,451) | 68.97% (5,390) |
| 10 Coronary atherosclerosis and other heart disease | Heart (Coronary arteries) | - | ✓ | 31.00% (2,452) | 69.00% (5,389) |
| 11 Diabetes mellitus with complications | Endocrine system (Pancreas) + Target organs | - | ✓ | 11.82% (909) | 88.18% (6,932) |
| 12 Diabetes mellitus without complication | Endocrine system (Pancreas) | - | ✓ | 20.58% (1,634) | 79.42% (6,207) |
| 13 Disorders of lipid metabolism | Liver and Circulatory system | - | ✓ | 41.05% (3,207) | 58.95% (4,634) |
| 14 Essential hypertension | Cardiovascular system | - | ✓ | 44.14% (3,467) | 55.86% (4,374) |
| 15 Fluid and electrolyte disorders | Kidneys and Endocrine system | - | ✓ | 45.48% (3,547) | 54.52% (4,294) |
| 16 Gastrointestinal hemorrhage | GI Tract | - | ✓ | 7.08% (562) | 92.92% (7,279) |
| 17 Hypertension with complications | Cardiovascular system | - | ✓ | 21.17% (1,649) | 78.83% (6,192) |
| 18 Other liver diseases | Liver | ✓ | ✓ | 16.12% (1,248) | 83.88% (6,593) |
| 19 Other lower respiratory disease | Lungs (Lower airways) | ✓ | ✓ | 12.93% (1,001) | 87.07% (6,840) |
| 20 Other upper respiratory disease | Lungs/Nasal passages (Upper airways) | ✓ | ✓ | 6.29% (500) | 93.71% (7,341) |
| 21 Pleurisy; pneumothorax | Lungs/Pleura | ✓ | ✓ | 9.93% (784) | 90.07% (7,057) |
| 22 Pneumonia (except caused by tuberculosis or std) | Lungs (Alveoli) | ✓ | ✓ | 18.84% (1,489) | 81.16% (6,352) |
| 23 Respiratory failure | Lungs and/or Neuromuscular system | - | ✓ | 28.37% (2,229) | 71.63% (5,612) |
| 24 Septicemia (except in labor) | Bloodstream (Systemic) | ✓ | ✓ | 22.21% (1,732) | 77.79% (6,109) |
| 25 Shock | Cardiovascular system (Systemic) | - | ✓ | 18.01% (1,413) | 81.99% (6,428) |

## 2. Training Details

**Hyperparameter Optimization and Model Selection**

We conducted 50 experiments corresponding to a randomly sampled hyperparameter set. Every set was evaluated across five runs with different random seeds. The learning rate was sampled uniformly from $[10^{-5}, 10^{-2}]$, and the number of training epochs from $\{5, 10, 15, 20, 30\}$. We used a fixed batch size of 16 and a cosine decay scheduler with $\alpha = 0$. Model checkpoints were saved at the final training epoch, and the optimal hyperparameter set was selected based on the highest mean validation AUROC across seeds, yielding a total of 250 runs per task.

Using the selected configuration, we retrained the models by combining the training and validation splits and ran five seeds with different initializations to obtain the final models for test evaluation. We report the performance results on the test set with means and standard errors computed over five random seeds. This procedure was applied to unimodal CXR, EHR, MedFuse, and MedFuse + GAP models following the original MedFuse protocol to identify the best-performing configuration based on AUROC.

All models, were trained using the Adam optimizer. Our experiments were executed using NVIDIA A100 and V100 80Gb Tensor Core GPUs.

**Data Availability**

The MIMIC (Medical Information Mart for Intensive Care) dataset is publicly available for research purposes. Access to the data requires completion of a data use agreement, which ensures compliance with the Health Insurance Portability and Accountability Act (HIPAA). Researchers can request access through the PhysioNet platform at https://physionet.org, where detailed instructions and requirements are provided. In our study, we used versions MIMIC-IV and MIMIC-CXR which include de-identified health data, vital signs, laboratory test results, medication records and chest X-ray images from patients admitted to the intensive care units at the Beth Israel Deaconess Medical Center.

## 3. Preliminaries

**Calibration.** A core requirement for selective prediction is that model confidence scores are well-calibrated, (i.e., predicted probabilities reflect true likelihoods of correctness), which is commonly measured using the Expected Calibration Error (ECE) (Guo et al., 2017). Formally, we can measure miscalibration as the expected gap between predicted confidence and empirical accuracy:

$$\text{ECE} = \mathbb{E}_{\hat{P}}\big[\,\big|\mathbb{P}(\hat{Y} = Y \mid \hat{P} = p) - p\big|\,\big]. \tag{3.1}$$

To approximate ECE in practice, we can use a finite partition of the probability space into $M$ bins $\{B_m\}_{m=1}^{M}$:

$$\widehat{\text{ECE}} = \sum_{m=1}^{M} \frac{|B_m|}{n} \big|\text{acc}(B_m) - \text{conf}(B_m)\big|, \tag{3.2}$$

where $\text{acc}(B_m)$ and $\text{conf}(B_m)$ denote the average accuracy and confidence of samples in bin $B_m$, respectively, and $n$ is the number of samples. Additionally, we may wish to compute the conditional calibration of a predictor. To do so, we can define a conditional expected calibration error,

$$\text{ECE}_c = \mathbb{E}_{\hat{P}|\hat{Y}=c}\big[\,\big|\mathbb{P}(\hat{Y} = c \mid \hat{P} = p) - p\big|\,\big]. \tag{3.3}$$

which captures class-dependent calibration. We can approximate the conditional calibration error by only considering the subset of points with label $c$ and, again, specifying a finite partition of the probability space into $M$ bins $\{B_m^c\}_{m=1}^{M}$:

$$\widehat{\text{ECE}}_c = \sum_{m=1}^{M} \frac{|B_m^c|}{n_c} \big|\text{acc}(B_m^c) - \text{conf}(B_m^c)\big|, \tag{3.4}$$

where $\text{acc}(B_m^c)$ and $\text{conf}(B_m^c)$ denote the average accuracy and confidence of samples in bin $B_m^c$, respectively, and $n_c$ is the number of samples whose label is $c$.

**Selective Prediction.** Selective prediction modifies the standard prediction pipeline by introducing a "reject option" denoted by $\perp$, using a gating mechanism defined by a selection function $s : \mathcal{X} \to \mathbb{R}$ that determines whether a prediction should be made for a given input point $\mathbf{x} \in \mathcal{X}$ (El-Yaniv et al., 2010). For a rejection threshold $\tau$, with $s$ representing the entropy of $\mathbf{x}$, the prediction model is given by:

$$(p(\mathbf{y} \mid \cdot; f), s)(\mathbf{x}) = \begin{cases} p(\mathbf{y} \mid \mathbf{x}; f), & \text{if } s \leq \tau \\ \perp, & \text{otherwise,} \end{cases} \tag{3.5}$$

with $p(\mathbf{y} \mid \cdot; f)$ being a model's predictive distribution. A variety of methods have been proposed to find a selection function $s$ (Rabanser et al., 2022; Rabanser and Papernot, 2025), with model's predictive uncertainty being a common choice: $s = \mathcal{H}(p(\mathbf{y} \mid \mathbf{x}; f))$. That is, a point $\mathbf{x} \in \mathcal{X}$ will be placed into the rejection class if the model's predictive uncertainty is above a certain threshold.

To evaluate the predictive performance of a prediction model $(p(\mathbf{y} \mid \cdot; f), s)(\mathbf{x})$, we compute the predictive performance of the classifier $p(\mathbf{y} \mid \mathbf{x}; f)$ over a range of thresholds $\tau$ (Figure 1.a). Performance is then assessed with standard metrics (e.g., accuracy, AUROC, AUPRC, or ECE), yielding selective prediction curves that capture the trade-off between coverage—the proportion of cases for which predictions are made—and performance (Geifman and El-Yaniv, 2017).

**Predictive Uncertainty Quantification.** Furthermore, understanding calibration requires examining how models quantify predictive uncertainty, which is typically decomposed into aleatoric uncertainty, arising from inherent noise in the data, and epistemic uncertainty, arising from limited knowledge about the model or parameters (Depeweg et al., 2018). A common way to quantify predictive uncertainty is via the Shannon entropy of the predictive distribution:

$$\mathcal{H}(p(\mathbf{y} \mid \mathbf{x})) = -\sum_c p(\mathbf{y} = c \mid \mathbf{x}) \log p(\mathbf{y} = c \mid \mathbf{x}), \tag{3.6}$$

where higher entropy indicates greater uncertainty (Rudner et al., 2023a,b).

Table 3: **Condition-level Performance of Standard and Selective Prediction Metrics**. MedFuse generally improves AUROC, AUPRC and selective AUROC, but calibration gains are inconsistent compared to EHR. (Dark-bold: $p < 0.05$, Wilcoxon signed-rank test, 5 seeds; Light-bold: highest mean, not significant)

| Clinical Condition | EHR | | | | | CXR | | | | | MedFuse | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | AUROC | AUPRC | Selective AUROC | Selective AUPRC | $\widehat{ECE}\downarrow$ | AUROC | AUPRC | Selective AUROC | Selective AUPRC | $\widehat{ECE}\downarrow$ | AUROC | AUPRC | Selective AUROC | Selective AUPRC | $\widehat{ECE}\downarrow$ |
| 1 | $0.724_{\pm0.003}$ | $0.564_{\pm0.010}$ | $0.772_{\pm0.100}$ | $0.627_{\pm0.100}$ | $\mathbf{1.62}_{\pm0.52}$ | $0.639_{\pm0.009}$ | $0.453_{\pm0.014}$ | $0.684_{\pm0.081}$ | $0.503_{\pm0.075}$ | $12.77_{\pm0.98}$ | $\mathbf{0.761}_{\pm0.005}$ | $\mathbf{0.589}_{\pm0.008}$ | $\mathbf{0.804}_{\pm0.111}$ | $\mathbf{0.634}_{\pm0.109}$ | $2.31_{\pm0.97}$ |
| 2 | $\mathbf{0.873}_{\pm0.005}$ | $\mathbf{0.459}_{\pm0.010}$ | $0.569_{\pm0.180}$ | $0.085_{\pm0.138}$ | $1.08_{\pm0.14}$ | $0.650_{\pm0.009}$ | $0.139_{\pm0.009}$ | $\mathbf{0.577}_{\pm0.121}$ | $0.084_{\pm0.057}$ | $3.99_{\pm0.31}$ | $0.856_{\pm0.010}$ | $0.432_{\pm0.019}$ | $0.560_{\pm0.176}$ | $0.075_{\pm0.116}$ | $\mathbf{0.88}_{\pm0.24}$ |
| 3 | $0.693_{\pm0.014}$ | $0.175_{\pm0.016}$ | $0.617_{\pm0.150}$ | $0.089_{\pm0.070}$ | $\mathbf{0.73}_{\pm0.07}$ | $0.640_{\pm0.015}$ | $0.153_{\pm0.017}$ | $0.567_{\pm0.102}$ | $\mathbf{0.092}_{\pm0.037}$ | $4.35_{\pm0.42}$ | $\mathbf{0.713}_{\pm0.005}$ | $\mathbf{0.209}_{\pm0.010}$ | $\mathbf{0.633}_{\pm0.207}$ | $0.083_{\pm0.049}$ | $0.76_{\pm0.23}$ |
| 4 | $0.598_{\pm0.004}$ | $0.462_{\pm0.012}$ | $0.618_{\pm0.060}$ | $0.491_{\pm0.065}$ | $\mathbf{1.15}_{\pm0.67}$ | $0.647_{\pm0.009}$ | $0.534_{\pm0.013}$ | $0.683_{\pm0.077}$ | $0.598_{\pm0.077}$ | $12.50_{\pm0.48}$ | $\mathbf{0.694}_{\pm0.006}$ | $\mathbf{0.565}_{\pm0.010}$ | $\mathbf{0.727}_{\pm0.079}$ | $\mathbf{0.610}_{\pm0.071}$ | $2.61_{\pm1.65}$ |
| 5 | $0.696_{\pm0.006}$ | $0.402_{\pm0.012}$ | $0.642_{\pm0.113}$ | $0.244_{\pm0.113}$ | $2.13_{\pm0.45}$ | $0.662_{\pm0.021}$ | $0.378_{\pm0.019}$ | $0.696_{\pm0.085}$ | $0.409_{\pm0.078}$ | $15.37_{\pm0.37}$ | $\mathbf{0.770}_{\pm0.003}$ | $\mathbf{0.519}_{\pm0.010}$ | $\mathbf{0.787}_{\pm0.177}$ | $\mathbf{0.521}_{\pm0.169}$ | $2.34_{\pm2.04}$ |
| 6 | $0.612_{\pm0.005}$ | $0.203_{\pm0.008}$ | $0.564_{\pm0.089}$ | $0.140_{\pm0.045}$ | $1.87_{\pm0.23}$ | $0.702_{\pm0.012}$ | $0.301_{\pm0.016}$ | $0.676_{\pm0.130}$ | $0.255_{\pm0.139}$ | $7.15_{\pm1.26}$ | $\mathbf{0.742}_{\pm0.004}$ | $\mathbf{0.345}_{\pm0.011}$ | $\mathbf{0.717}_{\pm0.120}$ | $\mathbf{0.279}_{\pm0.127}$ | $2.17_{\pm0.45}$ |
| 7 | $0.649_{\pm0.007}$ | $0.365_{\pm0.013}$ | $0.578_{\pm0.098}$ | $0.221_{\pm0.083}$ | $1.16_{\pm0.45}$ | $0.597_{\pm0.015}$ | $0.308_{\pm0.009}$ | $0.555_{\pm0.069}$ | $\mathbf{0.263}_{\pm0.067}$ | $9.15_{\pm0.58}$ | $\mathbf{0.673}_{\pm0.001}$ | $\mathbf{0.383}_{\pm0.003}$ | $\mathbf{0.602}_{\pm0.094}$ | $0.241_{\pm0.100}$ | $1.46_{\pm0.31}$ |
| 8 | $0.627_{\pm0.010}$ | $0.174_{\pm0.007}$ | $0.558_{\pm0.098}$ | $0.109_{\pm0.054}$ | $1.09_{\pm0.19}$ | $0.814_{\pm0.010}$ | $0.608_{\pm0.013}$ | $0.784_{\pm0.130}$ | $0.597_{\pm0.198}$ | $3.42_{\pm0.45}$ | $\mathbf{0.824}_{\pm0.004}$ | $\mathbf{0.609}_{\pm0.009}$ | $\mathbf{0.826}_{\pm0.099}$ | $\mathbf{0.648}_{\pm0.119}$ | $2.32_{\pm0.45}$ |
| 9 | $0.686_{\pm0.005}$ | $0.463_{\pm0.011}$ | $0.690_{\pm0.113}$ | $0.419_{\pm0.128}$ | $1.81_{\pm0.56}$ | $0.755_{\pm0.008}$ | $0.565_{\pm0.017}$ | $0.808_{\pm0.097}$ | $0.658_{\pm0.100}$ | $11.92_{\pm1.95}$ | $\mathbf{0.792}_{\pm0.004}$ | $\mathbf{0.625}_{\pm0.008}$ | $\mathbf{0.842}_{\pm0.110}$ | $\mathbf{0.703}_{\pm0.123}$ | $4.50_{\pm1.71}$ |
| 10 | $0.697_{\pm0.004}$ | $0.517_{\pm0.008}$ | $0.658_{\pm0.096}$ | $0.412_{\pm0.144}$ | $1.80_{\pm0.43}$ | $0.728_{\pm0.013}$ | $0.583_{\pm0.019}$ | $0.780_{\pm0.094}$ | $0.655_{\pm0.093}$ | $11.15_{\pm1.04}$ | $\mathbf{0.776}_{\pm0.002}$ | $\mathbf{0.651}_{\pm0.004}$ | $\mathbf{0.835}_{\pm0.109}$ | $\mathbf{0.723}_{\pm0.115}$ | $3.52_{\pm1.29}$ |
| 11 | $0.755_{\pm0.015}$ | $0.304_{\pm0.014}$ | $0.623_{\pm0.181}$ | $0.100_{\pm0.070}$ | $1.41_{\pm0.26}$ | $0.608_{\pm0.019}$ | $0.171_{\pm0.013}$ | $0.583_{\pm0.104}$ | $0.133_{\pm0.064}$ | $5.42_{\pm0.83}$ | $\mathbf{0.855}_{\pm0.003}$ | $\mathbf{0.485}_{\pm0.012}$ | $\mathbf{0.676}_{\pm0.269}$ | $\mathbf{0.198}_{\pm0.193}$ | $\mathbf{1.15}_{\pm0.22}$ |
| 12 | $0.623_{\pm0.012}$ | $0.286_{\pm0.009}$ | $0.556_{\pm0.085}$ | $0.195_{\pm0.061}$ | $1.38_{\pm0.28}$ | $0.589_{\pm0.022}$ | $0.257_{\pm0.013}$ | $0.590_{\pm0.088}$ | $\mathbf{0.214}_{\pm0.049}$ | $8.45_{\pm1.35}$ | $\mathbf{0.711}_{\pm0.006}$ | $\mathbf{0.355}_{\pm0.010}$ | $\mathbf{0.650}_{\pm0.121}$ | $0.200_{\pm0.110}$ | $1.74_{\pm0.56}$ |
| 13 | $0.650_{\pm0.003}$ | $0.553_{\pm0.010}$ | $0.642_{\pm0.085}$ | $0.620_{\pm0.070}$ | $\mathbf{0.99}_{\pm0.42}$ | $0.589_{\pm0.003}$ | $0.519_{\pm0.010}$ | $0.670_{\pm0.088}$ | $0.561_{\pm0.075}$ | $15.71_{\pm0.34}$ | $\mathbf{0.696}_{\pm0.002}$ | $\mathbf{0.599}_{\pm0.004}$ | $\mathbf{0.766}_{\pm0.139}$ | $\mathbf{0.668}_{\pm0.125}$ | $3.61_{\pm0.64}$ |
| 14 | $0.591_{\pm0.011}$ | $0.529_{\pm0.014}$ | $0.572_{\pm0.067}$ | $0.594_{\pm0.068}$ | $2.77_{\pm0.48}$ | $0.611_{\pm0.011}$ | $0.529_{\pm0.009}$ | $0.637_{\pm0.075}$ | $0.565_{\pm0.061}$ | $18.90_{\pm1.04}$ | $\mathbf{0.676}_{\pm0.004}$ | $\mathbf{0.603}_{\pm0.006}$ | $\mathbf{0.714}_{\pm0.087}$ | $\mathbf{0.660}_{\pm0.085}$ | $\mathbf{1.86}_{\pm0.62}$ |
| 15 | $0.703_{\pm0.005}$ | $0.676_{\pm0.005}$ | $0.715_{\pm0.103}$ | $0.782_{\pm0.083}$ | $1.41_{\pm0.42}$ | $0.600_{\pm0.011}$ | $0.544_{\pm0.009}$ | $0.626_{\pm0.071}$ | $0.586_{\pm0.062}$ | $16.93_{\pm1.17}$ | $\mathbf{0.718}_{\pm0.003}$ | $\mathbf{0.685}_{\pm0.004}$ | $\mathbf{0.755}_{\pm0.114}$ | $\mathbf{0.789}_{\pm0.082}$ | $3.39_{\pm0.92}$ |
| 16 | $0.628_{\pm0.007}$ | $0.111_{\pm0.008}$ | $0.590_{\pm0.144}$ | $0.073_{\pm0.033}$ | $\mathbf{0.58}_{\pm0.14}$ | $0.597_{\pm0.021}$ | $0.109_{\pm0.010}$ | $0.558_{\pm0.113}$ | $\mathbf{0.077}_{\pm0.028}$ | $4.36_{\pm0.27}$ | $\mathbf{0.665}_{\pm0.009}$ | $\mathbf{0.128}_{\pm0.003}$ | $\mathbf{0.606}_{\pm0.174}$ | $0.072_{\pm0.043}$ | $1.05_{\pm0.36}$ |
| 17 | $0.690_{\pm0.005}$ | $0.377_{\pm0.012}$ | $0.618_{\pm0.108}$ | $0.204_{\pm0.093}$ | $1.91_{\pm0.27}$ | $0.666_{\pm0.017}$ | $0.343_{\pm0.017}$ | $0.674_{\pm0.074}$ | $0.324_{\pm0.056}$ | $14.12_{\pm0.46}$ | $\mathbf{0.763}_{\pm0.004}$ | $\mathbf{0.482}_{\pm0.014}$ | $\mathbf{0.746}_{\pm0.188}$ | $\mathbf{0.413}_{\pm0.175}$ | $2.70_{\pm1.70}$ |
| 18 | $0.643_{\pm0.008}$ | $0.285_{\pm0.005}$ | $0.580_{\pm0.083}$ | $0.158_{\pm0.075}$ | $\mathbf{0.98}_{\pm0.19}$ | $0.672_{\pm0.013}$ | $0.345_{\pm0.003}$ | $\mathbf{0.674}_{\pm0.081}$ | $\mathbf{0.387}_{\pm0.112}$ | $6.42_{\pm0.82}$ | $\mathbf{0.712}_{\pm0.003}$ | $\mathbf{0.356}_{\pm0.009}$ | $0.658_{\pm0.101}$ | $0.273_{\pm0.097}$ | $1.75_{\pm0.19}$ |
| 19 | $0.594_{\pm0.008}$ | $0.172_{\pm0.008}$ | $0.543_{\pm0.080}$ | $0.120_{\pm0.032}$ | $\mathbf{0.55}_{\pm0.15}$ | $0.547_{\pm0.022}$ | $0.154_{\pm0.007}$ | $0.520_{\pm0.075}$ | $\mathbf{0.150}_{\pm0.065}$ | $5.55_{\pm1.19}$ | $\mathbf{0.611}_{\pm0.005}$ | $\mathbf{0.199}_{\pm0.008}$ | $\mathbf{0.568}_{\pm0.095}$ | $0.119_{\pm0.036}$ | $1.50_{\pm0.29}$ |
| 20 | $0.623_{\pm0.010}$ | $0.091_{\pm0.007}$ | $0.531_{\pm0.113}$ | $0.056_{\pm0.031}$ | $0.98_{\pm0.15}$ | $0.615_{\pm0.009}$ | $\mathbf{0.103}_{\pm0.009}$ | $0.523_{\pm0.091}$ | $\mathbf{0.060}_{\pm0.029}$ | $3.03_{\pm0.34}$ | $\mathbf{0.649}_{\pm0.011}$ | $0.098_{\pm0.003}$ | $\mathbf{0.532}_{\pm0.109}$ | $0.054_{\pm0.027}$ | $\mathbf{0.68}_{\pm0.10}$ |
| 21 | $0.591_{\pm0.009}$ | $0.129_{\pm0.006}$ | $0.573_{\pm0.129}$ | $0.097_{\pm0.037}$ | $1.91_{\pm0.38}$ | $0.658_{\pm0.023}$ | $\mathbf{0.179}_{\pm0.013}$ | $0.569_{\pm0.102}$ | $\mathbf{0.104}_{\pm0.048}$ | $5.11_{\pm0.98}$ | $\mathbf{0.663}_{\pm0.014}$ | $0.172_{\pm0.006}$ | $\mathbf{0.574}_{\pm0.103}$ | $0.089_{\pm0.056}$ | $1.43_{\pm0.48}$ |
| 22 | $0.722_{\pm0.007}$ | $0.368_{\pm0.007}$ | $0.607_{\pm0.119}$ | $0.184_{\pm0.114}$ | $\mathbf{1.10}_{\pm0.27}$ | $0.665_{\pm0.009}$ | $0.288_{\pm0.008}$ | $0.624_{\pm0.118}$ | $0.203_{\pm0.070}$ | $8.30_{\pm1.29}$ | $\mathbf{0.750}_{\pm0.004}$ | $\mathbf{0.402}_{\pm0.005}$ | $\mathbf{0.669}_{\pm0.154}$ | $\mathbf{0.214}_{\pm0.140}$ | $2.04_{\pm0.83}$ |
| 23 | $0.803_{\pm0.009}$ | $0.620_{\pm0.013}$ | $0.804_{\pm0.146}$ | $0.607_{\pm0.177}$ | $\mathbf{1.57}_{\pm0.36}$ | $0.714_{\pm0.008}$ | $0.493_{\pm0.017}$ | $0.725_{\pm0.102}$ | $0.475_{\pm0.095}$ | $11.41_{\pm1.37}$ | $\mathbf{0.817}_{\pm0.005}$ | $\mathbf{0.645}_{\pm0.007}$ | $\mathbf{0.872}_{\pm0.101}$ | $\mathbf{0.738}_{\pm0.115}$ | $2.09_{\pm0.85}$ |
| 24 | $0.754_{\pm0.006}$ | $0.483_{\pm0.010}$ | $0.691_{\pm0.141}$ | $0.358_{\pm0.177}$ | $1.88_{\pm0.40}$ | $0.649_{\pm0.005}$ | $0.350_{\pm0.008}$ | $0.626_{\pm0.099}$ | $0.302_{\pm0.097}$ | $11.07_{\pm1.03}$ | $\mathbf{0.776}_{\pm0.004}$ | $\mathbf{0.525}_{\pm0.010}$ | $\mathbf{0.752}_{\pm0.151}$ | $\mathbf{0.475}_{\pm0.178}$ | $2.38_{\pm0.68}$ |
| 25 | $0.808_{\pm0.006}$ | $0.519_{\pm0.023}$ | $0.704_{\pm0.179}$ | $0.316_{\pm0.214}$ | $1.55_{\pm0.48}$ | $0.693_{\pm0.005}$ | $0.336_{\pm0.006}$ | $0.659_{\pm0.110}$ | $0.301_{\pm0.106}$ | $8.71_{\pm1.31}$ | $\mathbf{0.820}_{\pm0.003}$ | $\mathbf{0.554}_{\pm0.005}$ | $\mathbf{0.748}_{\pm0.203}$ | $\mathbf{0.413}_{\pm0.242}$ | $1.65_{\pm0.35}$ |
| Average | $0.681_{\pm0.007}$ | $0.371_{\pm0.010}$ | $0.625_{\pm0.115}$ | $0.292_{\pm0.091}$ | $\mathbf{1.42}_{\pm0.34}$ | $0.654_{\pm0.013}$ | $0.350_{\pm0.012}$ | $0.643_{\pm0.095}$ | $0.342_{\pm0.077}$ | $9.41_{\pm0.83}$ | $\mathbf{0.739}_{\pm0.005}$ | $\mathbf{0.449}_{\pm0.008}$ | $\mathbf{0.705}_{\pm0.136}$ | $\mathbf{0.396}_{\pm0.112}$ | $2.08_{\pm0.73}$ |

## 4. Further Empirical Results

In this section, we present complimentary findings of our study on selective prediction and calibration in multimodal clinical condition classification guided by the following main questions:

1. *Do multimodal models improve or degrade reliable selective prediction?*
2. *Do other multimodal architectures exhibit the same reliability patterns?*
3. *Can loss-upweight training mitigate calibration failures?*

To this end, we compare unimodal and multimodal models on standard discrimination and calibration metrics to establish whether multimodal fusion improves reliability. Additionally, we decompose calibration into positive- and negative-class components to identify which types of predictions drive miscalibration. We also contextualize these calibration properties with selective prediction behavior across coverage thresholds, showing when and why selective performance degrades. This analysis allows us to assess both aggregate and fine-grained metrics, linking calibration quality with selective reliability. We conclude by evaluating a simple label-upweighted loss training strategy, which provides improvements in selective prediction but highlights the remaining challenges in mitigating class-dependent miscalibration.

### 4.1. *Do Multimodal Models Improve or Degrade Reliable Selective Prediction?*

**Overall Performance.** Table 3 presents the performance metrics for the unimodal models and the deterministic multimodal baseline (MedFuse). Across the 25 conditions, MedFuse generally achieves higher AUROC, AUPRC, and selective AUROC than both unimodal baselines, often with statistical significance. However, its calibration gains over EHR are less consistent, illustrating a recurrent mismatch between discrimination and probability alignment.

For example, in *Acute and unspecified renal failure (1)*, MedFuse achieves the highest AUROC (0.761) and AUPRC (0.589), reducing ECE relative to CXR (12.77 vs. 2.31), yet EHR remains the best calibrated

Table 4: **Condition-level Performance of ECE stratification.** Positive-class ECE dominates across conditions, highlighting overconfidence in positive predictions. (Dark-bold: $p < 0.05$, Wilcoxon signed-rank test, 5 seeds; Light-bold: highest mean, not significant)

| Clinical Condition | EHR | | | CXR | | | MedFuse | | |
|---|---|---|---|---|---|---|---|---|---|
| | $\widehat{\mathrm{ECE}} \downarrow$ | $\widehat{\mathrm{ECE}}_{c=1} \downarrow$ | $\widehat{\mathrm{ECE}}_{c=0} \downarrow$ | $\widehat{\mathrm{ECE}} \downarrow$ | $\widehat{\mathrm{ECE}}_{c=1} \downarrow$ | $\widehat{\mathrm{ECE}}_{c=0} \downarrow$ | $\widehat{\mathrm{ECE}} \downarrow$ | $\widehat{\mathrm{ECE}}_{c=1} \downarrow$ | $\widehat{\mathrm{ECE}}_{c=0} \downarrow$ |
| 1 | **1.62**$_{\pm 0.52}$ | 30.76$_{\pm 1.82}$ | 14.32$_{\pm 1.14}$ | 12.77$_{\pm 0.98}$ | 35.26$_{\pm 5.07}$ | **3.63**$_{\pm 1.74}$ | 2.31$_{\pm 0.97}$ | **24.85**$_{\pm 3.96}$ | 9.97$_{\pm 1.14}$ |
| 2 | 1.08$_{\pm 0.14}$ | **47.76**$_{\pm 3.00}$ | 4.68$_{\pm 0.25}$ | 3.99$_{\pm 0.31}$ | 83.21$_{\pm 3.28}$ | 4.13$_{\pm 0.75}$ | **0.88**$_{\pm 0.24}$ | 50.56$_{\pm 3.07}$ | **3.56**$_{\pm 0.19}$ |
| 3 | **0.73**$_{\pm 0.07}$ | 88.27$_{\pm 0.33}$ | 8.22$_{\pm 0.06}$ | 4.35$_{\pm 0.42}$ | 80.59$_{\pm 5.85}$ | 4.91$_{\pm 0.70}$ | 0.76$_{\pm 0.23}$ | 86.93$_{\pm 0.14}$ | 8.05$_{\pm 0.25}$ |
| 4 | **1.15**$_{\pm 0.67}$ | 49.64$_{\pm 1.14}$ | 28.93$_{\pm 0.65}$ | 12.50$_{\pm 0.48}$ | 29.02$_{\pm 5.95}$ | **4.34**$_{\pm 3.01}$ | 2.61$_{\pm 1.65}$ | **26.90**$_{\pm 9.52}$ | 13.57$_{\pm 2.63}$ |
| 5 | **2.13**$_{\pm 0.45}$ | 67.12$_{\pm 0.76}$ | 19.71$_{\pm 0.41}$ | 15.37$_{\pm 0.37}$ | 45.19$_{\pm 2.51}$ | 5.92$_{\pm 0.92}$ | 2.34$_{\pm 2.04}$ | **40.71**$_{\pm 8.56}$ | 10.21$_{\pm 1.06}$ |
| 6 | **1.87**$_{\pm 0.23}$ | 81.50$_{\pm 0.22}$ | 16.00$_{\pm 0.07}$ | 7.15$_{\pm 1.26}$ | 64.69$_{\pm 14.12}$ | 3.32$_{\pm 1.04}$ | 2.17$_{\pm 0.45}$ | **61.34**$_{\pm 3.97}$ | 9.27$_{\pm 0.68}$ |
| 7 | **1.16**$_{\pm 0.45}$ | 70.71$_{\pm 1.80}$ | 20.18$_{\pm 0.20}$ | 9.15$_{\pm 0.58}$ | 66.07$_{\pm 3.97}$ | 7.68$_{\pm 1.09}$ | 1.46$_{\pm 0.31}$ | 66.13$_{\pm 1.23}$ | 19.12$_{\pm 0.40}$ |
| 8 | **1.09**$_{\pm 0.19}$ | 87.20$_{\pm 0.28}$ | 10.42$_{\pm 0.04}$ | 3.42$_{\pm 0.45}$ | 38.99$_{\pm 6.15}$ | 2.77$_{\pm 0.79}$ | 2.32$_{\pm 0.45}$ | 41.64$_{\pm 0.88}$ | 4.38$_{\pm 0.89}$ |
| 9 | **1.81**$_{\pm 0.56}$ | 41.15$_{\pm 1.10}$ | 17.36$_{\pm 0.89}$ | 11.92$_{\pm 1.95}$ | 20.38$_{\pm 5.59}$ | 8.40$_{\pm 4.88}$ | 4.50$_{\pm 1.71}$ | 24.13$_{\pm 9.07}$ | **4.69**$_{\pm 1.96}$ |
| 10 | **1.80**$_{\pm 0.43}$ | 38.09$_{\pm 0.54}$ | 17.50$_{\pm 0.76}$ | 11.15$_{\pm 1.04}$ | 30.34$_{\pm 5.39}$ | 3.44$_{\pm 2.02}$ | 3.52$_{\pm 1.29}$ | **22.44**$_{\pm 4.96}$ | 7.08$_{\pm 1.88}$ |
| 11 | 1.41$_{\pm 0.26}$ | 80.37$_{\pm 1.11}$ | 10.92$_{\pm 0.15}$ | 5.42$_{\pm 0.83}$ | 74.92$_{\pm 6.19}$ | 5.25$_{\pm 0.72}$ | **1.15**$_{\pm 0.22}$ | **47.13**$_{\pm 2.20}$ | 6.30$_{\pm 0.34}$ |
| 12 | **1.38**$_{\pm 0.28}$ | 77.14$_{\pm 0.38}$ | 19.78$_{\pm 0.07}$ | 8.45$_{\pm 1.35}$ | 62.23$_{\pm 7.64}$ | 6.34$_{\pm 2.18}$ | 1.74$_{\pm 0.56}$ | 68.84$_{\pm 2.85}$ | 17.64$_{\pm 0.79}$ |
| 13 | **0.99**$_{\pm 0.42}$ | 22.03$_{\pm 1.04}$ | 15.95$_{\pm 0.73}$ | 15.71$_{\pm 0.34}$ | 12.50$_{\pm 9.56}$ | 18.16$_{\pm 6.99}$ | 3.61$_{\pm 0.64}$ | 12.51$_{\pm 4.83}$ | **4.88**$_{\pm 0.85}$ |
| 14 | 2.77$_{\pm 0.48}$ | 26.80$_{\pm 2.53}$ | 25.28$_{\pm 1.99}$ | 18.90$_{\pm 1.04}$ | 24.94$_{\pm 8.12}$ | 14.45$_{\pm 7.33}$ | **1.86**$_{\pm 0.62}$ | **12.64**$_{\pm 5.92}$ | 9.09$_{\pm 4.73}$ |
| 15 | **1.41**$_{\pm 0.42}$ | **10.36**$_{\pm 0.73}$ | 8.89$_{\pm 0.74}$ | 16.93$_{\pm 1.17}$ | 26.86$_{\pm 3.30}$ | 8.69$_{\pm 1.27}$ | 3.39$_{\pm 0.92}$ | 12.05$_{\pm 2.66}$ | **3.95**$_{\pm 0.52}$ |
| 16 | **0.58**$_{\pm 0.14}$ | 91.64$_{\pm 0.15}$ | 6.89$_{\pm 0.09}$ | 4.36$_{\pm 0.27}$ | 87.04$_{\pm 3.51}$ | 5.26$_{\pm 0.63}$ | 1.05$_{\pm 0.36}$ | 89.64$_{\pm 0.80}$ | 7.21$_{\pm 0.42}$ |
| 17 | **1.91**$_{\pm 0.27}$ | 72.75$_{\pm 0.85}$ | 18.83$_{\pm 0.48}$ | 14.12$_{\pm 0.48}$ | 50.40$_{\pm 1.45}$ | 4.05$_{\pm 0.56}$ | 2.70$_{\pm 1.70}$ | **46.22**$_{\pm 7.91}$ | 10.70$_{\pm 0.96}$ |
| 18 | **0.98**$_{\pm 0.19}$ | 80.92$_{\pm 0.15}$ | 15.41$_{\pm 0.14}$ | 6.42$_{\pm 0.82}$ | 59.33$_{\pm 5.58}$ | 5.30$_{\pm 0.43}$ | 1.75$_{\pm 0.19}$ | 60.84$_{\pm 4.55}$ | 11.05$_{\pm 0.92}$ |
| 19 | **0.55**$_{\pm 0.15}$ | 86.08$_{\pm 0.27}$ | 12.65$_{\pm 0.12}$ | 5.55$_{\pm 1.19}$ | 78.99$_{\pm 4.96}$ | 7.68$_{\pm 0.63}$ | 1.50$_{\pm 0.29}$ | 86.71$_{\pm 0.43}$ | 11.15$_{\pm 0.47}$ |
| 20 | 0.98$_{\pm 0.15}$ | 91.94$_{\pm 0.23}$ | 6.23$_{\pm 0.09}$ | 3.03$_{\pm 0.34}$ | 86.07$_{\pm 1.83}$ | 3.03$_{\pm 0.70}$ | **0.68**$_{\pm 0.10}$ | 92.06$_{\pm 0.39}$ | 5.53$_{\pm 0.33}$ |
| 21 | 1.91$_{\pm 0.38}$ | 88.61$_{\pm 0.18}$ | 9.83$_{\pm 0.14}$ | 5.11$_{\pm 0.98}$ | 81.78$_{\pm 4.69}$ | 3.49$_{\pm 0.52}$ | **1.43**$_{\pm 0.48}$ | 87.93$_{\pm 0.48}$ | 8.56$_{\pm 0.60}$ |
| 22 | **1.10**$_{\pm 0.27}$ | 63.95$_{\pm 0.75}$ | 14.87$_{\pm 0.30}$ | 8.30$_{\pm 1.29}$ | 64.87$_{\pm 8.07}$ | 4.72$_{\pm 0.81}$ | 2.04$_{\pm 0.83}$ | **61.13**$_{\pm 2.02}$ | 11.72$_{\pm 0.84}$ |
| 23 | **1.57**$_{\pm 0.36}$ | **24.30**$_{\pm 1.74}$ | 10.96$_{\pm 0.99}$ | 11.41$_{\pm 1.37}$ | 43.39$_{\pm 7.84}$ | **2.39**$_{\pm 0.73}$ | 2.09$_{\pm 0.85}$ | 25.38$_{\pm 3.17}$ | 7.17$_{\pm 0.15}$ |
| 24 | **1.88**$_{\pm 0.40}$ | 47.22$_{\pm 1.96}$ | 13.14$_{\pm 0.80}$ | 11.07$_{\pm 1.03}$ | 59.06$_{\pm 5.83}$ | 3.53$_{\pm 1.01}$ | 2.38$_{\pm 0.68}$ | **43.16**$_{\pm 1.67}$ | 9.60$_{\pm 0.84}$ |
| 25 | **1.55**$_{\pm 0.48}$ | 42.32$_{\pm 3.98}$ | 9.78$_{\pm 0.65}$ | 8.71$_{\pm 0.31}$ | 60.56$_{\pm 4.27}$ | **2.58**$_{\pm 1.21}$ | 1.65$_{\pm 0.35}$ | **38.14**$_{\pm 1.48}$ | 7.12$_{\pm 0.45}$ |
| Average | **1.42**$_{\pm 0.34}$ | 60.35$_{\pm 1.08}$ | 14.27$_{\pm 0.48}$ | 9.41$_{\pm 0.83}$ | 54.67$_{\pm 5.63}$ | **5.74**$_{\pm 1.71}$ | 2.08$_{\pm 0.73}$ | **49.20**$_{\pm 3.47}$ | 8.86$_{\pm 0.97}$ |

model (1.62). In *Coronary atherosclerosis (10)*, MedFuse nearly doubles EHR's calibration error (3.52 vs. 1.80) but presents the best performance across all other metrics. In fact, the only condition where MedFuse outperforms both unimodal variants in calibration at a statistically significant level is *Other upper respiratory disease (20)* (MedFuse = 0.68 vs. CXR = 3.03, EHR = 0.98). Yet in this case, the improvements do not translate into statistically superior AUROC, AUPRC, or their selective prediction counterparts.

Overall, multimodal fusion reliably boosts discrimination but fails to improve, and sometimes worsens, calibration. This highlights the risk of relying solely on standard metrics when assessing multimodal models in safety-critical settings.

**Class-level Calibration Stratification.** To probe this behavior, we decompose ECE into positive- and negative-class components as defined by Equation (3.3) (i.e., $\mathrm{ECE}_{c=1}$ and $\mathrm{ECE}_{c=0}$). As shown in Table 4, we observe that positive-class ECE is consistently larger than negative-class ECE, showing that overconfidence in positive predictions is the dominant calibration failure. This imbalance is hidden when reporting only standard ECE and can therefore mislead conclusions about robustness.

We also observe that no single model dominates across all conditions. EHR and MedFuse generally achieve lower positive-class ECE than CXR, while CXR often achieves the lowest negative-class ECE. This inversion illustrates how aggregate calibration masks strengths and weaknesses in specific components and why positive-class miscalibration, in particular, undermines selective prediction. It also suggests that the multimodal improvements in positive-class calibration likely originate from the EHR modality rather than from fusion per se.

Table 5: **Condition-level Performance of Multimodal Architectures**. Although the performance of all three models is comparable, MedFuse reliably outperforms both DrFuse and MeTra across different conditions and evaluation metrics, with the downside that all models present high ECE scores. This further confirms that complex architectures do not necessarily improve predictive performance and do not solve the reliability issues observed in this clinical task. (Dark-bold: $p < 0.05$, Wilcoxon signed-rank test, 5 seeds; Light-bold: highest mean, not significant)

| Clinical Condition | MedFuse | | | | | DrFuse | | | | | MeTra | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | AUROC | AUPRC | Selective AUROC | Selective AUPRC | $\widehat{\text{ECE}}\downarrow$ | AUROC | AUPRC | Selective AUROC | Selective AUPRC | $\widehat{\text{ECE}}\downarrow$ | AUROC | AUPRC | Selective AUROC | Selective AUPRC | $\widehat{\text{ECE}}\downarrow$ |
| 1 | $0.761_{\pm0.005}$ | $0.589_{\pm0.008}$ | $0.804_{\pm0.111}$ | $0.634_{\pm0.109}$ | $2.31_{\pm0.97}$ | $0.732_{\pm0.004}$ | $0.551_{\pm0.008}$ | $0.767_{\pm0.104}$ | $0.585_{\pm0.113}$ | $4.06_{\pm1.21}$ | $0.722_{\pm0.009}$ | $0.541_{\pm0.011}$ | $0.759_{\pm0.083}$ | $0.573_{\pm0.071}$ | $2.73_{\pm0.39}$ |
| 2 | $0.856_{\pm0.010}$ | $0.432_{\pm0.019}$ | $0.560_{\pm0.176}$ | $0.075_{\pm0.116}$ | $0.88_{\pm0.24}$ | $0.899_{\pm0.004}$ | $0.460_{\pm0.013}$ | $0.584_{\pm0.272}$ | $0.118_{\pm0.166}$ | $1.06_{\pm0.44}$ | $0.900_{\pm0.004}$ | $0.490_{\pm0.022}$ | $0.602_{\pm0.253}$ | $0.112_{\pm0.172}$ | $1.13_{\pm0.32}$ |
| 3 | $0.713_{\pm0.005}$ | $0.209_{\pm0.010}$ | $0.633_{\pm0.207}$ | $0.083_{\pm0.049}$ | $0.76_{\pm0.23}$ | $0.697_{\pm0.013}$ | $0.166_{\pm0.010}$ | $0.637_{\pm0.169}$ | $0.090_{\pm0.069}$ | $2.21_{\pm0.58}$ | $0.717_{\pm0.010}$ | $0.197_{\pm0.014}$ | $0.610_{\pm0.158}$ | $0.085_{\pm0.058}$ | $1.98_{\pm1.12}$ |
| 4 | $0.694_{\pm0.006}$ | $0.565_{\pm0.010}$ | $0.727_{\pm0.079}$ | $0.610_{\pm0.071}$ | $2.61_{\pm1.65}$ | $0.650_{\pm0.012}$ | $0.520_{\pm0.013}$ | $0.692_{\pm0.084}$ | $0.547_{\pm0.089}$ | $3.42_{\pm1.66}$ | $0.628_{\pm0.011}$ | $0.500_{\pm0.008}$ | $0.648_{\pm0.076}$ | $0.529_{\pm0.079}$ | $3.59_{\pm1.65}$ |
| 5 | $0.770_{\pm0.003}$ | $0.519_{\pm0.010}$ | $0.787_{\pm0.177}$ | $0.521_{\pm0.169}$ | $2.34_{\pm2.04}$ | $0.726_{\pm0.006}$ | $0.448_{\pm0.009}$ | $0.693_{\pm0.127}$ | $0.334_{\pm0.152}$ | $3.59_{\pm2.23}$ | $0.709_{\pm0.017}$ | $0.414_{\pm0.019}$ | $0.668_{\pm0.101}$ | $0.299_{\pm0.105}$ | $2.68_{\pm1.31}$ |
| 6 | $0.742_{\pm0.004}$ | $0.345_{\pm0.011}$ | $0.717_{\pm0.120}$ | $0.279_{\pm0.127}$ | $2.17_{\pm0.45}$ | $0.670_{\pm0.013}$ | $0.274_{\pm0.015}$ | $0.573_{\pm0.099}$ | $0.137_{\pm0.060}$ | $1.94_{\pm0.66}$ | $0.618_{\pm0.010}$ | $0.219_{\pm0.007}$ | $0.557_{\pm0.086}$ | $0.140_{\pm0.042}$ | $2.52_{\pm0.61}$ |
| 7 | $0.673_{\pm0.001}$ | $0.383_{\pm0.003}$ | $0.602_{\pm0.094}$ | $0.241_{\pm0.100}$ | $1.46_{\pm0.31}$ | $0.693_{\pm0.008}$ | $0.409_{\pm0.008}$ | $0.662_{\pm0.097}$ | $0.349_{\pm0.097}$ | $1.90_{\pm0.93}$ | $0.701_{\pm0.003}$ | $0.427_{\pm0.005}$ | $0.673_{\pm0.091}$ | $0.380_{\pm0.104}$ | $1.76_{\pm0.51}$ |
| 8 | $0.824_{\pm0.004}$ | $0.609_{\pm0.009}$ | $0.826_{\pm0.099}$ | $0.648_{\pm0.119}$ | $2.32_{\pm0.45}$ | $0.691_{\pm0.007}$ | $0.234_{\pm0.006}$ | $0.611_{\pm0.121}$ | $0.104_{\pm0.066}$ | $2.30_{\pm0.70}$ | $0.674_{\pm0.007}$ | $0.220_{\pm0.011}$ | $0.592_{\pm0.111}$ | $0.107_{\pm0.051}$ | $2.61_{\pm0.72}$ |
| 9 | $0.792_{\pm0.004}$ | $0.625_{\pm0.008}$ | $0.842_{\pm0.110}$ | $0.703_{\pm0.123}$ | $4.50_{\pm1.71}$ | $0.732_{\pm0.006}$ | $0.544_{\pm0.009}$ | $0.751_{\pm0.137}$ | $0.527_{\pm0.169}$ | $3.95_{\pm2.06}$ | $0.705_{\pm0.013}$ | $0.503_{\pm0.019}$ | $0.723_{\pm0.103}$ | $0.515_{\pm0.120}$ | $2.02_{\pm0.62}$ |
| 10 | $0.776_{\pm0.002}$ | $0.651_{\pm0.004}$ | $0.835_{\pm0.109}$ | $0.723_{\pm0.115}$ | $3.52_{\pm1.29}$ | $0.735_{\pm0.005}$ | $0.575_{\pm0.005}$ | $0.764_{\pm0.123}$ | $0.574_{\pm0.154}$ | $4.17_{\pm3.96}$ | $0.705_{\pm0.004}$ | $0.545_{\pm0.009}$ | $0.734_{\pm0.093}$ | $0.560_{\pm0.102}$ | $3.28_{\pm0.71}$ |
| 11 | $0.855_{\pm0.006}$ | $0.485_{\pm0.012}$ | $0.676_{\pm0.269}$ | $0.198_{\pm0.193}$ | $1.15_{\pm0.22}$ | $0.867_{\pm0.005}$ | $0.502_{\pm0.015}$ | $0.775_{\pm0.267}$ | $0.464_{\pm0.226}$ | $1.41_{\pm0.45}$ | $0.838_{\pm0.008}$ | $0.450_{\pm0.027}$ | $0.822_{\pm0.192}$ | $0.487_{\pm0.234}$ | $2.22_{\pm1.11}$ |
| 12 | $0.711_{\pm0.006}$ | $0.355_{\pm0.010}$ | $0.650_{\pm0.121}$ | $0.200_{\pm0.110}$ | $1.74_{\pm0.56}$ | $0.726_{\pm0.007}$ | $0.377_{\pm0.007}$ | $0.673_{\pm0.131}$ | $0.262_{\pm0.132}$ | $3.52_{\pm1.24}$ | $0.700_{\pm0.020}$ | $0.361_{\pm0.024}$ | $0.620_{\pm0.116}$ | $0.223_{\pm0.117}$ | $2.22_{\pm0.45}$ |
| 13 | $0.696_{\pm0.002}$ | $0.599_{\pm0.004}$ | $0.766_{\pm0.139}$ | $0.668_{\pm0.125}$ | $3.61_{\pm0.64}$ | $0.678_{\pm0.005}$ | $0.596_{\pm0.010}$ | $0.706_{\pm0.089}$ | $0.670_{\pm0.094}$ | $3.72_{\pm1.65}$ | $0.650_{\pm0.004}$ | $0.568_{\pm0.011}$ | $0.682_{\pm0.083}$ | $0.647_{\pm0.085}$ | $2.57_{\pm1.67}$ |
| 14 | $0.676_{\pm0.005}$ | $0.603_{\pm0.006}$ | $0.714_{\pm0.087}$ | $0.660_{\pm0.085}$ | $1.86_{\pm0.62}$ | $0.647_{\pm0.008}$ | $0.580_{\pm0.013}$ | $0.639_{\pm0.086}$ | $0.650_{\pm0.074}$ | $2.02_{\pm1.20}$ | $0.627_{\pm0.011}$ | $0.560_{\pm0.009}$ | $0.623_{\pm0.087}$ | $0.620_{\pm0.073}$ | $2.58_{\pm1.81}$ |
| 15 | $0.718_{\pm0.004}$ | $0.685_{\pm0.004}$ | $0.755_{\pm0.114}$ | $0.789_{\pm0.082}$ | $3.39_{\pm0.92}$ | $0.720_{\pm0.004}$ | $0.676_{\pm0.005}$ | $0.738_{\pm0.099}$ | $0.773_{\pm0.074}$ | $3.78_{\pm0.92}$ | $0.698_{\pm0.007}$ | $0.660_{\pm0.017}$ | $0.742_{\pm0.096}$ | $0.754_{\pm0.085}$ | $3.54_{\pm0.48}$ |
| 16 | $0.665_{\pm0.009}$ | $0.128_{\pm0.003}$ | $0.606_{\pm0.174}$ | $0.072_{\pm0.043}$ | $1.05_{\pm0.36}$ | $0.711_{\pm0.006}$ | $0.171_{\pm0.009}$ | $0.625_{\pm0.221}$ | $0.065_{\pm0.044}$ | $1.01_{\pm0.19}$ | $0.690_{\pm0.017}$ | $0.159_{\pm0.023}$ | $0.625_{\pm0.168}$ | $0.069_{\pm0.055}$ | $1.34_{\pm0.49}$ |
| 17 | $0.763_{\pm0.004}$ | $0.482_{\pm0.014}$ | $0.746_{\pm0.188}$ | $0.413_{\pm0.175}$ | $2.70_{\pm1.70}$ | $0.714_{\pm0.006}$ | $0.413_{\pm0.008}$ | $0.665_{\pm0.123}$ | $0.281_{\pm0.138}$ | $3.71_{\pm2.60}$ | $0.701_{\pm0.017}$ | $0.386_{\pm0.018}$ | $0.635_{\pm0.103}$ | $0.247_{\pm0.098}$ | $2.52_{\pm0.89}$ |
| 18 | $0.712_{\pm0.003}$ | $0.356_{\pm0.010}$ | $0.658_{\pm0.101}$ | $0.273_{\pm0.097}$ | $1.75_{\pm0.19}$ | $0.700_{\pm0.006}$ | $0.315_{\pm0.008}$ | $0.625_{\pm0.115}$ | $0.167_{\pm0.094}$ | $2.10_{\pm0.48}$ | $0.674_{\pm0.010}$ | $0.288_{\pm0.016}$ | $0.612_{\pm0.113}$ | $0.157_{\pm0.073}$ | $2.37_{\pm1.06}$ |
| 19 | $0.611_{\pm0.005}$ | $0.199_{\pm0.008}$ | $0.568_{\pm0.095}$ | $0.119_{\pm0.036}$ | $1.50_{\pm0.29}$ | $0.598_{\pm0.013}$ | $0.171_{\pm0.006}$ | $0.556_{\pm0.082}$ | $0.120_{\pm0.032}$ | $1.57_{\pm0.58}$ | $0.575_{\pm0.016}$ | $0.157_{\pm0.016}$ | $0.538_{\pm0.080}$ | $0.124_{\pm0.036}$ | $1.67_{\pm0.44}$ |
| 20 | $0.649_{\pm0.011}$ | $0.098_{\pm0.003}$ | $0.532_{\pm0.109}$ | $0.054_{\pm0.027}$ | $0.68_{\pm0.10}$ | $0.694_{\pm0.015}$ | $0.139_{\pm0.016}$ | $0.565_{\pm0.137}$ | $0.053_{\pm0.055}$ | $1.23_{\pm0.58}$ | $0.671_{\pm0.016}$ | $0.113_{\pm0.013}$ | $0.554_{\pm0.130}$ | $0.057_{\pm0.056}$ | $1.29_{\pm0.44}$ |
| 21 | $0.663_{\pm0.014}$ | $0.172_{\pm0.006}$ | $0.574_{\pm0.103}$ | $0.089_{\pm0.056}$ | $1.43_{\pm0.48}$ | $0.649_{\pm0.003}$ | $0.166_{\pm0.007}$ | $0.569_{\pm0.108}$ | $0.093_{\pm0.046}$ | $1.37_{\pm0.77}$ | $0.606_{\pm0.010}$ | $0.131_{\pm0.007}$ | $0.596_{\pm0.114}$ | $0.097_{\pm0.036}$ | $2.26_{\pm0.22}$ |
| 22 | $0.750_{\pm0.004}$ | $0.402_{\pm0.005}$ | $0.669_{\pm0.154}$ | $0.214_{\pm0.140}$ | $2.04_{\pm0.83}$ | $0.762_{\pm0.002}$ | $0.423_{\pm0.008}$ | $0.668_{\pm0.152}$ | $0.209_{\pm0.138}$ | $2.20_{\pm1.12}$ | $0.740_{\pm0.008}$ | $0.388_{\pm0.019}$ | $0.649_{\pm0.162}$ | $0.210_{\pm0.124}$ | $2.00_{\pm1.18}$ |
| 23 | $0.817_{\pm0.005}$ | $0.645_{\pm0.007}$ | $0.872_{\pm0.101}$ | $0.738_{\pm0.115}$ | $2.09_{\pm0.85}$ | $0.819_{\pm0.003}$ | $0.621_{\pm0.008}$ | $0.829_{\pm0.151}$ | $0.630_{\pm0.169}$ | $2.94_{\pm1.00}$ | $0.818_{\pm0.005}$ | $0.620_{\pm0.019}$ | $0.800_{\pm0.170}$ | $0.543_{\pm0.209}$ | $2.36_{\pm0.82}$ |
| 24 | $0.776_{\pm0.004}$ | $0.525_{\pm0.010}$ | $0.752_{\pm0.151}$ | $0.475_{\pm0.178}$ | $2.38_{\pm0.68}$ | $0.797_{\pm0.005}$ | $0.545_{\pm0.012}$ | $0.776_{\pm0.155}$ | $0.513_{\pm0.197}$ | $2.76_{\pm0.80}$ | $0.783_{\pm0.008}$ | $0.528_{\pm0.015}$ | $0.747_{\pm0.155}$ | $0.436_{\pm0.190}$ | $2.57_{\pm0.80}$ |
| 25 | $0.820_{\pm0.003}$ | $0.554_{\pm0.005}$ | $0.748_{\pm0.203}$ | $0.413_{\pm0.242}$ | $1.65_{\pm0.35}$ | $0.847_{\pm0.004}$ | $0.571_{\pm0.014}$ | $0.772_{\pm0.233}$ | $0.484_{\pm0.253}$ | $1.71_{\pm0.67}$ | $0.830_{\pm0.008}$ | $0.552_{\pm0.026}$ | $0.719_{\pm0.205}$ | $0.367_{\pm0.236}$ | $2.18_{\pm1.00}$ |
| Average | $0.739_{\pm0.005}$ | $0.449_{\pm0.008}$ | $0.705_{\pm0.136}$ | $0.396_{\pm0.112}$ | $2.08_{\pm0.73}$ | $0.726_{\pm0.007}$ | $0.418_{\pm0.010}$ | $0.677_{\pm0.139}$ | $0.352_{\pm0.116}$ | $2.54_{\pm1.15}$ | $0.707_{\pm0.010}$ | $0.399_{\pm0.015}$ | $0.661_{\pm0.125}$ | $0.333_{\pm0.104}$ | $2.32_{\pm0.83}$ |

**Takeaway 1.** Multimodal fusion improves discrimination but does not reliably improve calibration; positive-class overconfidence is the dominant source of error.

## 4.2. *Do Other Multimodal Architectures Exhibit the Same Reliability Patterns?*

**Extending to DrFuse and MeTRA.** To evaluate whether our findings generalize beyond MedFuse, we additionally assess two multimodal architectures, DrFuse (Yao et al., 2024) and MeTRA (Khader et al., 2023), which also integrate EHR and CXR modalities for clinical condition prediction. Despite distinct and more complex design choices, both models continue to exhibit reliability trends consistent with MedFuse. In particular, Table 5 shows that MedFuse is consistently the best performing model in terms of evaluation metrics, without any single architecture meaningfully being best in terms of calibration error across conditions.

Extending our correlation analysis between positive-class ECE AUC and Selective AUC, in Figure C1 we show that all three multimodal architectures exhibit the same behaviour for both AUROC and AUPRC, highlighting the strong positive-class ECE, maintaining the same imbalance of overconfidence in positive predictions affecting selective AUC.

This is further demonstrated by the ECE stratification results across models (Table 6), showing that complex fusion architectures do not tackle the reliability issues of class-dependent miscalibration. In addition, the selective prediction curves (Appendix 4.5, Figures C8 and C9) follow nearly identical patterns to MedFuse: moderate gains at intermediate rejection thresholds followed by collapse at extreme coverage levels, suggesting that the underlying calibration dynamics, not the specific fusion mechanism, drive the observed reliability outcomes. Finally, in Appendix 6 (Figure D1, Figure D2, Figure D3), we visualize three representative clinical conditions (2, 8, 25) comparing all multimodal models, confirming that while DrFuse and MeTRA

Table 6: **Condition-level Performance of ECE stratification Across Multimodal Architectures.** Positive-class ECE dominates across conditions, highlighting overconfidence for in positive predictions for all multimodal architectures regardless of the complexity of the fusion strategy. (Dark-bold: $p < 0.05$, Wilcoxon signed-rank test, 5 seeds; Light-bold: highest mean, not significant)

| Clinical Condition | MedFuse | | | DrFuse | | | MeTra | | |
|---|---|---|---|---|---|---|---|---|---|
| | $\widehat{\text{ECE}} \downarrow$ | $\widehat{\text{ECE}}_{c=1} \downarrow$ | $\widehat{\text{ECE}}_{c=0} \downarrow$ | $\widehat{\text{ECE}} \downarrow$ | $\widehat{\text{ECE}}_{c=1} \downarrow$ | $\widehat{\text{ECE}}_{c=0} \downarrow$ | $\widehat{\text{ECE}} \downarrow$ | $\widehat{\text{ECE}}_{c=1} \downarrow$ | $\widehat{\text{ECE}}_{c=0} \downarrow$ |
| 1 | $\mathbf{2.31}_{\pm 0.97}$ | $\mathbf{24.85}_{\pm 3.96}$ | $9.97_{\pm 1.14}$ | $4.06_{\pm 1.21}$ | $27.74_{\pm 5.78}$ | $\mathbf{7.01}_{\pm 2.40}$ | $2.73_{\pm 0.39}$ | $35.27_{\pm 2.04}$ | $13.65_{\pm 1.70}$ |
| 2 | $\mathbf{0.88}_{\pm 0.24}$ | $50.56_{\pm 3.07}$ | $3.56_{\pm 0.19}$ | $1.06_{\pm 0.44}$ | $37.08_{\pm 7.82}$ | $\mathbf{2.98}_{\pm 0.46}$ | $1.13_{\pm 0.32}$ | $\mathbf{33.53}_{\pm 10.49}$ | $3.52_{\pm 0.15}$ |
| 3 | $\mathbf{0.76}_{\pm 0.23}$ | $86.93_{\pm 0.14}$ | $8.05_{\pm 0.25}$ | $2.21_{\pm 0.58}$ | $88.15_{\pm 2.27}$ | $\mathbf{7.83}_{\pm 1.71}$ | $1.98_{\pm 1.12}$ | $\mathbf{84.65}_{\pm 2.99}$ | $9.32_{\pm 1.95}$ |
| 4 | $\mathbf{2.61}_{\pm 1.65}$ | $\mathbf{26.90}_{\pm 9.52}$ | $\mathbf{13.57}_{\pm 2.63}$ | $3.42_{\pm 1.66}$ | $35.79_{\pm 13.68}$ | $17.54_{\pm 7.67}$ | $3.59_{\pm 1.65}$ | $39.75_{\pm 13.97}$ | $19.53_{\pm 5.75}$ |
| 5 | $2.34_{\pm 2.04}$ | $\mathbf{40.71}_{\pm 8.56}$ | $\mathbf{10.21}_{\pm 1.06}$ | $3.59_{\pm 2.23}$ | $49.62_{\pm 9.38}$ | $12.31_{\pm 1.96}$ | $2.68_{\pm 1.31}$ | $60.85_{\pm 4.41}$ | $16.17_{\pm 0.66}$ |
| 6 | $2.17_{\pm 0.45}$ | $\mathbf{61.34}_{\pm 3.97}$ | $\mathbf{9.27}_{\pm 0.68}$ | $\mathbf{1.94}_{\pm 0.66}$ | $79.53_{\pm 1.91}$ | $13.44_{\pm 1.82}$ | $2.52_{\pm 0.61}$ | $80.43_{\pm 0.66}$ | $15.52_{\pm 1.17}$ |
| 7 | $1.46_{\pm 0.31}$ | $66.13_{\pm 1.23}$ | $19.12_{\pm 0.40}$ | $1.90_{\pm 0.93}$ | $65.82_{\pm 3.70}$ | $18.38_{\pm 1.68}$ | $1.76_{\pm 0.51}$ | $\mathbf{56.89}_{\pm 6.12}$ | $\mathbf{16.28}_{\pm 0.59}$ |
| 8 | $2.32_{\pm 0.45}$ | $\mathbf{41.64}_{\pm 0.88}$ | $\mathbf{4.38}_{\pm 0.89}$ | $\mathbf{2.30}_{\pm 0.70}$ | $82.47_{\pm 3.45}$ | $9.75_{\pm 1.62}$ | $2.61_{\pm 0.72}$ | $82.51_{\pm 4.15}$ | $10.35_{\pm 1.64}$ |
| 9 | $4.50_{\pm 1.71}$ | $24.13_{\pm 9.07}$ | $4.69_{\pm 1.96}$ | $3.95_{\pm 2.06}$ | $34.55_{\pm 11.33}$ | $10.07_{\pm 2.26}$ | $\mathbf{2.02}_{\pm 0.62}$ | $38.30_{\pm 9.77}$ | $14.65_{\pm 3.27}$ |
| 10 | $3.52_{\pm 1.29}$ | $22.44_{\pm 4.96}$ | $7.08_{\pm 1.88}$ | $4.17_{\pm 3.96}$ | $34.09_{\pm 13.14}$ | $11.16_{\pm 2.60}$ | $\mathbf{3.28}_{\pm 0.71}$ | $37.68_{\pm 7.23}$ | $16.13_{\pm 1.71}$ |
| 11 | $1.15_{\pm 0.22}$ | $47.13_{\pm 2.20}$ | $6.30_{\pm 0.34}$ | $1.41_{\pm 0.45}$ | $\mathbf{35.95}_{\pm 5.61}$ | $\mathbf{5.01}_{\pm 0.89}$ | $2.22_{\pm 1.11}$ | $55.17_{\pm 11.94}$ | $6.25_{\pm 1.24}$ |
| 12 | $\mathbf{1.74}_{\pm 0.56}$ | $68.84_{\pm 2.85}$ | $17.64_{\pm 0.79}$ | $3.52_{\pm 1.24}$ | $\mathbf{55.32}_{\pm 11.18}$ | $\mathbf{16.14}_{\pm 3.87}$ | $2.22_{\pm 0.45}$ | $62.38_{\pm 8.01}$ | $16.60_{\pm 1.96}$ |
| 13 | $3.61_{\pm 0.64}$ | $\mathbf{12.51}_{\pm 4.83}$ | $\mathbf{4.88}_{\pm 0.85}$ | $3.72_{\pm 1.65}$ | $17.27_{\pm 12.41}$ | $8.27_{\pm 5.89}$ | $\mathbf{2.57}_{\pm 1.67}$ | $27.10_{\pm 8.21}$ | $15.39_{\pm 2.24}$ |
| 14 | $1.86_{\pm 0.62}$ | $\mathbf{12.64}_{\pm 5.92}$ | $9.09_{\pm 4.73}$ | $2.02_{\pm 1.20}$ | $16.93_{\pm 11.53}$ | $12.38_{\pm 6.58}$ | $2.58_{\pm 1.81}$ | $14.50_{\pm 10.83}$ | $10.32_{\pm 7.64}$ |
| 15 | $3.39_{\pm 0.92}$ | $12.05_{\pm 2.66}$ | $3.95_{\pm 0.52}$ | $3.78_{\pm 0.92}$ | $\mathbf{10.59}_{\pm 3.32}$ | $\mathbf{3.01}_{\pm 1.87}$ | $3.54_{\pm 0.48}$ | $14.38_{\pm 6.59}$ | $7.02_{\pm 3.27}$ |
| 16 | $1.05_{\pm 0.36}$ | $89.64_{\pm 0.80}$ | $7.21_{\pm 0.42}$ | $\mathbf{1.01}_{\pm 0.19}$ | $88.91_{\pm 0.80}$ | $6.84_{\pm 0.74}$ | $1.34_{\pm 0.49}$ | $\mathbf{88.77}_{\pm 2.63}$ | $6.89_{\pm 0.98}$ |
| 17 | $2.70_{\pm 1.70}$ | $\mathbf{46.22}_{\pm 7.91}$ | $\mathbf{10.70}_{\pm 0.96}$ | $3.71_{\pm 2.60}$ | $54.74_{\pm 9.45}$ | $12.52_{\pm 2.57}$ | $\mathbf{2.52}_{\pm 0.89}$ | $65.36_{\pm 3.14}$ | $15.65_{\pm 0.56}$ |
| 18 | $\mathbf{1.75}_{\pm 0.19}$ | $60.84_{\pm 4.55}$ | $\mathbf{11.05}_{\pm 0.92}$ | $2.10_{\pm 0.48}$ | $70.47_{\pm 3.99}$ | $14.39_{\pm 1.54}$ | $2.37_{\pm 1.06}$ | $78.01_{\pm 3.50}$ | $14.10_{\pm 1.39}$ |
| 19 | $1.50_{\pm 0.29}$ | $86.71_{\pm 0.43}$ | $11.15_{\pm 0.47}$ | $1.57_{\pm 0.58}$ | $86.99_{\pm 1.08}$ | $11.54_{\pm 1.14}$ | $1.67_{\pm 0.44}$ | $\mathbf{86.28}_{\pm 0.82}$ | $12.37_{\pm 0.55}$ |
| 20 | $0.68_{\pm 0.10}$ | $92.06_{\pm 0.39}$ | $5.53_{\pm 0.33}$ | $1.23_{\pm 0.58}$ | $91.04_{\pm 0.93}$ | $\mathbf{4.63}_{\pm 0.81}$ | $1.29_{\pm 0.44}$ | $\mathbf{90.78}_{\pm 2.42}$ | $5.28_{\pm 1.12}$ |
| 21 | $1.43_{\pm 0.48}$ | $\mathbf{87.93}_{\pm 0.48}$ | $8.56_{\pm 0.60}$ | $\mathbf{1.37}_{\pm 0.77}$ | $89.23_{\pm 0.93}$ | $8.20_{\pm 0.81}$ | $2.26_{\pm 0.22}$ | $88.25_{\pm 1.49}$ | $9.52_{\pm 1.33}$ |
| 22 | $2.04_{\pm 0.83}$ | $61.13_{\pm 2.02}$ | $11.72_{\pm 0.84}$ | $2.20_{\pm 1.12}$ | $\mathbf{59.74}_{\pm 7.21}$ | $\mathbf{11.26}_{\pm 1.04}$ | $\mathbf{2.00}_{\pm 1.18}$ | $62.64_{\pm 6.49}$ | $12.51_{\pm 1.24}$ |
| 23 | $\mathbf{2.09}_{\pm 0.85}$ | $25.38_{\pm 3.17}$ | $7.17_{\pm 0.15}$ | $2.94_{\pm 1.00}$ | $22.12_{\pm 3.08}$ | $4.87_{\pm 0.58}$ | $2.36_{\pm 0.82}$ | $\mathbf{18.30}_{\pm 5.56}$ | $5.82_{\pm 1.39}$ |
| 24 | $2.38_{\pm 0.68}$ | $43.16_{\pm 1.67}$ | $9.60_{\pm 0.84}$ | $2.76_{\pm 0.80}$ | $36.35_{\pm 4.61}$ | $7.43_{\pm 1.94}$ | $2.57_{\pm 0.80}$ | $45.99_{\pm 7.53}$ | $10.29_{\pm 1.44}$ |
| 25 | $\mathbf{1.65}_{\pm 0.35}$ | $38.14_{\pm 1.48}$ | $7.12_{\pm 0.45}$ | $1.71_{\pm 0.67}$ | $\mathbf{28.90}_{\pm 2.32}$ | $\mathbf{4.83}_{\pm 1.20}$ | $2.18_{\pm 1.00}$ | $40.13_{\pm 7.60}$ | $7.04_{\pm 0.57}$ |
| **Average** | $2.08_{\pm 0.73}$ | $49.20_{\pm 3.47}$ | $8.86_{\pm 0.97}$ | $2.54_{\pm 1.15}$ | $51.98_{\pm 6.04}$ | $9.67_{\pm 2.23}$ | $2.32_{\pm 0.83}$ | $55.52_{\pm 5.94}$ | $11.61_{\pm 1.82}$ |

occasionally show local improvements, the overall calibration and selective prediction trends remain virtually unchanged.

> **Takeaway 2.** Advances in multimodal fusion—such as transformer-based models—do not alter the fundamental trend: multimodal architectures enhance discrimination but leave class-dependent miscalibration largely unresolved.

### 4.3. *Can Loss-Upweight Training Mitigate Calibration Failures?*

**Label-Dependent Loss Upweighting.** To probe and mitigate class-dependent miscalibration, we implement a simple label-dependent loss reweighting procedure. For each clinical condition $c$, we assign a weight $w_c$ to its loss term, increasing the contribution of low-prevalence (underrepresented) positive labels during training, which explicitly upweights rare but clinically important conditions. For each multimodal backbone (MedFuse, DrFuse, MeTra), we train a label-upweighted version, enabling us to study how this simple intervention affects class-dependent calibration and selective prediction across architectures and capacities.

**Persistent Reliability Gaps.** To test whether the observed miscalibration can be mitigated, we utilize a straightforward method using positive prevalence as a scaling loss factor during training to emphasize the model's attention on underrepresented conditions. Table 7 shows the improvement of selective prediction

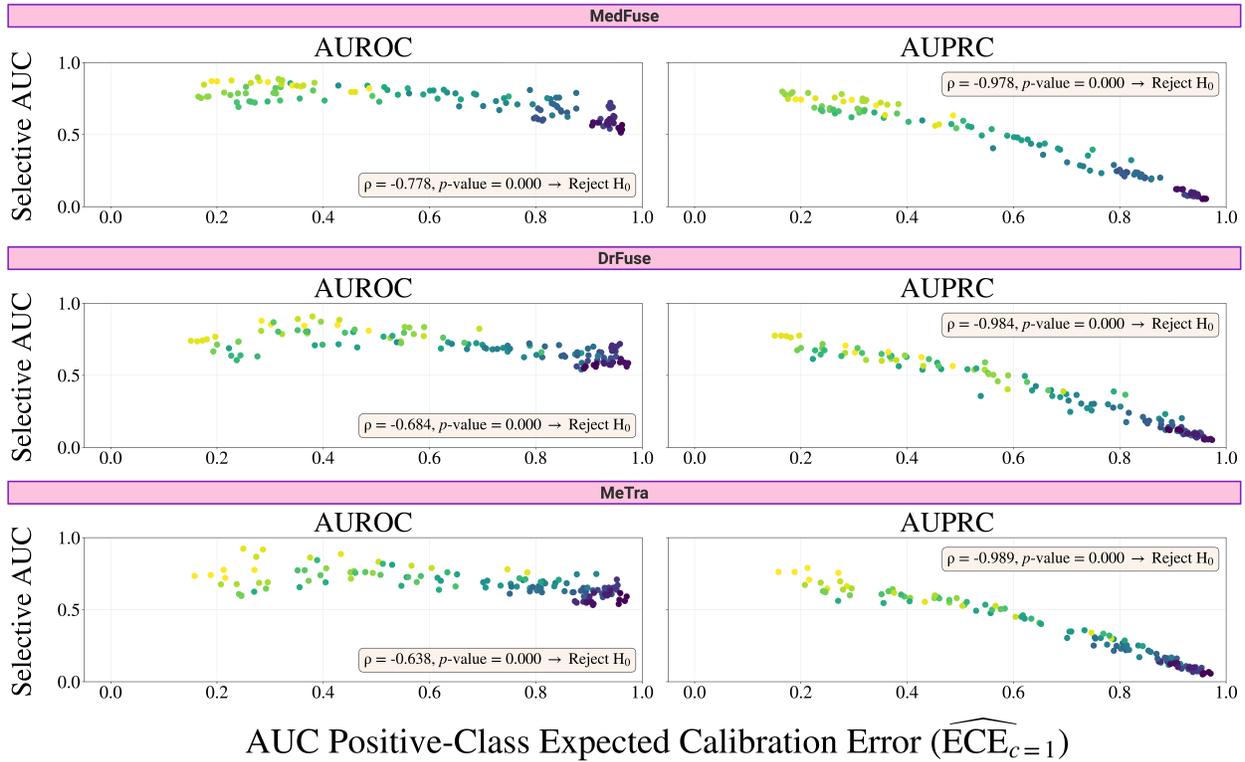AUC Positive-Class Expected Calibration Error ($\widehat{\text{ECE}}_{c=1}$)

Figure C1: **Do Other Multimodal Architectures Exhibit the Same Reliability Patterns?** For each multimodal architecture, higher positive-class ECE AUC consistently predicts lower selective AUC across conditions for both AUROC and AUPRC, indicating that overconfidence in positive predictions drives instability in selective evaluation regardless of fusion strategy. For all architectures and metrics, we reject the null hypothesis of no non-negative association between stratified ECE and selective AUC (Spearman's rank $H_0$, $p < 0.05$), showing that increased architectural complexity alone does not resolve performance–reliability trade-offs.

metrics over the standard AUC when using loss upweighting for the MedFuse architecture. We note consistent but modest improvements in selective AUROC and AUPRC, with $\Delta$AUC (Selective–Standard). While average $\Delta$AUPRC values are slightly higher, the improvements remain small and variable across conditions. Still, these improvements fall short of statistical significance, indicating that while selective behaviour trends upward, the underlying reliability gap persists.

Turning to calibration, standard ECE computation shows no global improvement, with only minor and statistically insignificant increases in error for the upweighted variant. However, stratified ECE analysis paints a more nuanced picture: the upweighted model achieves statistically significant positive-class calibration improvements in 23 of 25 conditions, and negative-class improvements in 11 of 25. These findings confirm that loss reweighting effectively reduces overconfidence for rare positive labels, though this does not fully translate into significant enhancements in selective prediction.

Furthermore, when applying the same loss upweight strategy to DrFuse and MeTra we observe similar patterns. In particular, Figure C2 shows that the negative relationship between positive-class ECE AUC and selective AUC has some general improvements as points are clustered closer together in lower ECE AUC values. However, we consistently reject our defined $H_0$ (no non-negative correlation) with statistical significance which shows that although there are some improvements, this simple loss mechanism does not fully resolve the reliability issues in selective prediction for clinical condition classification.

Table 7: **Loss Upweighting During Training Leads to Inconsistent Improvements.** Improvements are concentrated mainly in stratified ECE but the reliability issue in improvement between standard and selective AUC remains. (Dark-bold: $p < 0.05$, Wilcoxon signed-rank test, 5 seeds)

| Clinical Condition | MedFuse Upscaled | | Improvement (MedFuse vs MedFuse Upscaled) | | | | |
|---|---|---|---|---|---|---|---|
| | Δ AUROC | Δ AUPRC | Δ Selective AUROC | Δ Selective AUPRC | $\Delta\widehat{ECE}\downarrow$ | $\Delta\widehat{ECE}_{c=1}\downarrow$ | $\Delta\widehat{ECE}_{c=0}\downarrow$ |
| 1 | $-0.004_{\pm0.016}$ | $\mathbf{0.084}_{\pm0.006}$ | $-0.047_{\pm0.031}$ | $0.038_{\pm0.064}$ | $0.008_{\pm0.013}$ | $\mathbf{-0.210}_{\pm0.041}$ | $\mathbf{-0.050}_{\pm0.033}$ |
| 2 | $0.001_{\pm0.008}$ | $\mathbf{0.059}_{\pm0.006}$ | $\mathbf{0.312}_{\pm0.015}$ | $\mathbf{0.415}_{\pm0.018}$ | $0.015_{\pm0.010}$ | $\mathbf{-0.397}_{\pm0.035}$ | $\mathbf{-0.015}_{\pm0.008}$ |
| 3 | $\mathbf{0.028}_{\pm0.009}$ | $\mathbf{0.054}_{\pm0.015}$ | $\mathbf{0.111}_{\pm0.005}$ | $\mathbf{0.176}_{\pm0.018}$ | $0.088_{\pm0.022}$ | $\mathbf{-0.748}_{\pm0.023}$ | $0.033_{\pm0.028}$ |
| 4 | $-0.042_{\pm0.025}$ | $\mathbf{0.068}_{\pm0.005}$ | $-0.077_{\pm0.019}$ | $0.021_{\pm0.032}$ | $0.003_{\pm0.005}$ | $\mathbf{-0.194}_{\pm0.090}$ | $\mathbf{-0.082}_{\pm0.046}$ |
| 5 | $\mathbf{0.045}_{\pm0.018}$ | $\mathbf{0.103}_{\pm0.008}$ | $0.018_{\pm0.048}$ | $0.078_{\pm0.101}$ | $0.023_{\pm0.017}$ | $\mathbf{-0.357}_{\pm0.099}$ | $-0.048_{\pm0.044}$ |
| 6 | $\mathbf{0.055}_{\pm0.033}$ | $\mathbf{0.097}_{\pm0.010}$ | $0.085_{\pm0.062}$ | $\mathbf{0.165}_{\pm0.068}$ | $0.026_{\pm0.027}$ | $\mathbf{-0.475}_{\pm0.055}$ | $-0.033_{\pm0.042}$ |
| 7 | $-0.025_{\pm0.003}$ | $\mathbf{0.085}_{\pm0.011}$ | $\mathbf{0.043}_{\pm0.009}$ | $\mathbf{0.235}_{\pm0.026}$ | $0.044_{\pm0.016}$ | $\mathbf{-0.563}_{\pm0.026}$ | $\mathbf{-0.087}_{\pm0.017}$ |
| 8 | $\mathbf{0.041}_{\pm0.012}$ | $\mathbf{0.150}_{\pm0.016}$ | $0.034_{\pm0.032}$ | $0.086_{\pm0.103}$ | $0.023_{\pm0.008}$ | $\mathbf{-0.215}_{\pm0.059}$ | $0.020_{\pm0.013}$ |
| 9 | $\mathbf{0.051}_{\pm0.028}$ | $\mathbf{0.112}_{\pm0.012}$ | $-0.002_{\pm0.012}$ | $0.023_{\pm0.039}$ | $0.015_{\pm0.012}$ | $\mathbf{-0.175}_{\pm0.103}$ | $0.032_{\pm0.058}$ |
| 10 | $0.036_{\pm0.031}$ | $\mathbf{0.091}_{\pm0.010}$ | $-0.026_{\pm0.023}$ | $0.009_{\pm0.028}$ | $0.023_{\pm0.009}$ | $\mathbf{-0.171}_{\pm0.075}$ | $0.004_{\pm0.047}$ |
| 11 | $\mathbf{0.040}_{\pm0.014}$ | $\mathbf{0.124}_{\pm0.027}$ | $\mathbf{0.225}_{\pm0.029}$ | $\mathbf{0.410}_{\pm0.072}$ | $0.035_{\pm0.008}$ | $\mathbf{-0.423}_{\pm0.036}$ | $\mathbf{-0.013}_{\pm0.006}$ |
| 12 | $0.007_{\pm0.013}$ | $\mathbf{0.046}_{\pm0.018}$ | $\mathbf{0.070}_{\pm0.015}$ | $\mathbf{0.212}_{\pm0.027}$ | $0.051_{\pm0.009}$ | $\mathbf{-0.607}_{\pm0.048}$ | $\mathbf{-0.072}_{\pm0.032}$ |
| 13 | $\mathbf{0.044}_{\pm0.030}$ | $\mathbf{0.065}_{\pm0.014}$ | $-0.032_{\pm0.028}$ | $-0.013_{\pm0.018}$ | $0.023_{\pm0.007}$ | $-0.053_{\pm0.058}$ | $0.054_{\pm0.044}$ |
| 14 | $0.018_{\pm0.036}$ | $\mathbf{0.059}_{\pm0.008}$ | $-0.030_{\pm0.017}$ | $-0.001_{\pm0.019}$ | $0.004_{\pm0.010}$ | $-0.048_{\pm0.083}$ | $-0.026_{\pm0.085}$ |
| 15 | $0.012_{\pm0.014}$ | $\mathbf{0.096}_{\pm0.004}$ | $-0.028_{\pm0.022}$ | $-0.017_{\pm0.010}$ | $-0.002_{\pm0.014}$ | $\mathbf{-0.078}_{\pm0.022}$ | $0.004_{\pm0.022}$ |
| 16 | $\mathbf{0.076}_{\pm0.012}$ | $\mathbf{0.042}_{\pm0.004}$ | $\mathbf{0.125}_{\pm0.006}$ | $\mathbf{0.098}_{\pm0.006}$ | $0.062_{\pm0.041}$ | $\mathbf{-0.769}_{\pm0.068}$ | $0.004_{\pm0.041}$ |
| 17 | $\mathbf{0.051}_{\pm0.017}$ | $\mathbf{0.102}_{\pm0.006}$ | $0.057_{\pm0.049}$ | $0.150_{\pm0.110}$ | $0.024_{\pm0.014}$ | $\mathbf{-0.401}_{\pm0.099}$ | $-0.048_{\pm0.044}$ |
| 18 | $0.005_{\pm0.026}$ | $\mathbf{0.101}_{\pm0.017}$ | $0.065_{\pm0.052}$ | $\mathbf{0.195}_{\pm0.080}$ | $0.040_{\pm0.030}$ | $\mathbf{-0.522}_{\pm0.060}$ | $-0.042_{\pm0.044}$ |
| 19 | $0.011_{\pm0.011}$ | $\mathbf{0.052}_{\pm0.006}$ | $\mathbf{0.055}_{\pm0.008}$ | $\mathbf{0.132}_{\pm0.006}$ | $0.022_{\pm0.011}$ | $\mathbf{-0.807}_{\pm0.032}$ | $\mathbf{-0.073}_{\pm0.019}$ |
| 20 | $0.008_{\pm0.009}$ | $\mathbf{0.056}_{\pm0.008}$ | $\mathbf{0.175}_{\pm0.016}$ | $\mathbf{0.146}_{\pm0.006}$ | $0.038_{\pm0.017}$ | $\mathbf{-0.739}_{\pm0.016}$ | $-0.004_{\pm0.021}$ |
| 21 | $-0.046_{\pm0.031}$ | $\mathbf{0.027}_{\pm0.012}$ | $\mathbf{0.079}_{\pm0.049}$ | $\mathbf{0.116}_{\pm0.026}$ | $0.051_{\pm0.038}$ | $\mathbf{-0.804}_{\pm0.041}$ | $-0.012_{\pm0.044}$ |
| 22 | $\mathbf{0.025}_{\pm0.016}$ | $\mathbf{0.086}_{\pm0.020}$ | $\mathbf{0.120}_{\pm0.012}$ | $\mathbf{0.287}_{\pm0.025}$ | $0.019_{\pm0.003}$ | $\mathbf{-0.559}_{\pm0.020}$ | $\mathbf{-0.062}_{\pm0.026}$ |
| 23 | $0.007_{\pm0.012}$ | $\mathbf{0.104}_{\pm0.004}$ | $-0.047_{\pm0.010}$ | $0.019_{\pm0.029}$ | $0.001_{\pm0.010}$ | $\mathbf{-0.193}_{\pm0.043}$ | $\mathbf{-0.039}_{\pm0.020}$ |
| 24 | $\mathbf{0.041}_{\pm0.012}$ | $\mathbf{0.110}_{\pm0.003}$ | $\mathbf{0.063}_{\pm0.033}$ | $\mathbf{0.146}_{\pm0.050}$ | $0.032_{\pm0.017}$ | $\mathbf{-0.379}_{\pm0.026}$ | $\mathbf{-0.029}_{\pm0.015}$ |
| 25 | $\mathbf{0.060}_{\pm0.008}$ | $\mathbf{0.142}_{\pm0.012}$ | $\mathbf{0.138}_{\pm0.021}$ | $\mathbf{0.285}_{\pm0.061}$ | $0.034_{\pm0.007}$ | $\mathbf{-0.325}_{\pm0.021}$ | $\mathbf{-0.020}_{\pm0.012}$ |

**Takeaway 3.** Loss upweighting for underrepresented labels yields measurable calibration gains, particularly for positive predictions, but fails to produce consistent or statistically significant improvements in selective reliability.
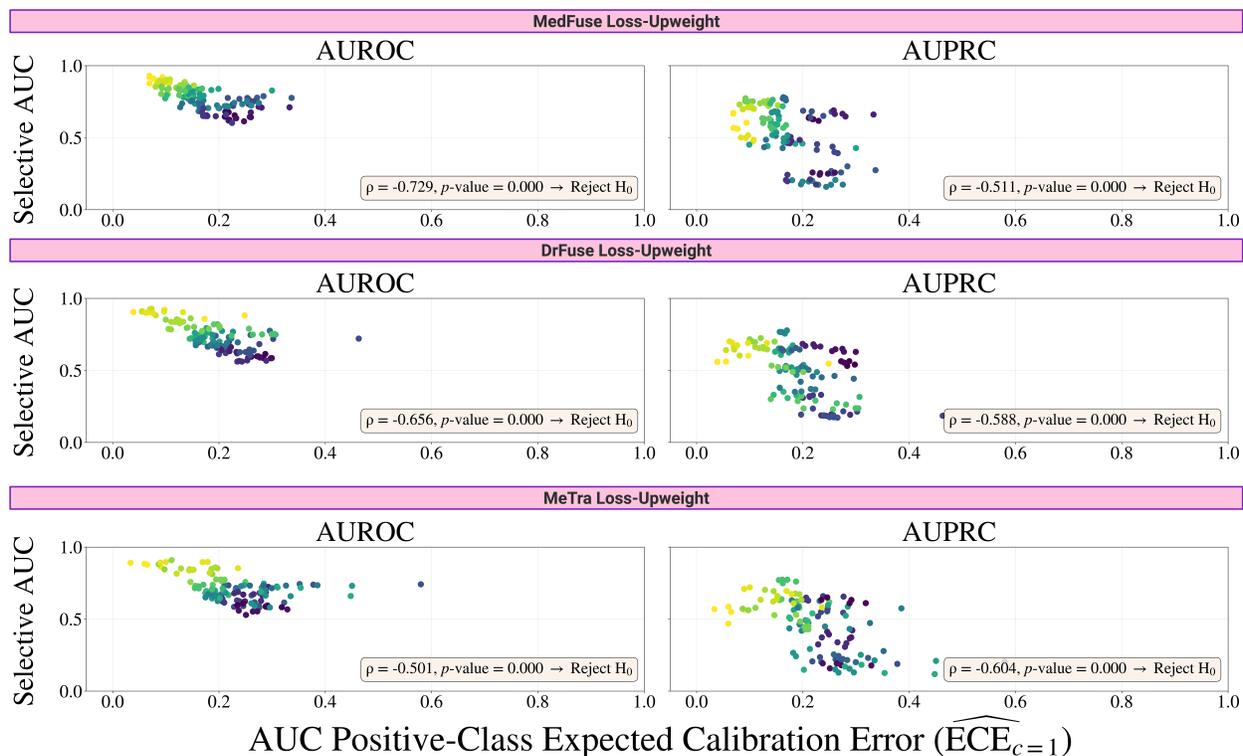
Figure C2: **Can Loss-Upweight Training Mitigate Calibration Failures?** For the loss-upweighted versions of MedFuse, DrFuse, and MeTra, the negative relationship between positive-class ECE AUC and selective AUC remains: conditions with higher stratified ECE still exhibit lower selective AUROC and AUPRC. While points cluster toward lower ECE AUC values, indicating partial calibration improvements, we consistently reject the null hypothesis of no non-negative association between ECE and selective AUC (Spearman's rank $p < 0.05$), showing that this simple loss reweighting does not fully resolve reliability failures in selective prediction for any of the models.

## 4.4. Regression Analysis of Calibration and Discriminative Performance

To strengthen the link between calibration and selective performance, we fit linear regression models on MedFuse predictions across all 25 clinical conditions using class-stratified ECE (positive and negative) as predictors of both AUROC and AUPRC. Results reveal that stratified ECE, particularly the positive-class component, exhibits a strong negative correlation with AUPRC, reinforcing that miscalibrated confidence in positive predictions reliably forecasts degraded precision and recall under selective evaluation.



Figure C3: **Linear Regression of AUROC Against Stratified ECE.** Relationship between AUROC and both positive- and negative-class ECE across 25 conditions. Regression slopes demonstrate a mild but consistent negative association, indicating that poorer calibration modestly reduces discriminative performance.

Figure C4: **Linear Regression of AUPRC Against Stratified ECE.** Relationship between AUPRC and class-stratified ECE across 25 conditions. The correlation between ECE and AUPRC exhibits a clear, monotonic decline—confirming that high positive-class miscalibration strongly predicts lower precision-recall performance under selective prediction.

## 4.5. Stratified Calibration Behavior Across Modalities

To complement the results in Figure 2, we visualize class-stratified ECE and selective AUC for unimodal (CXR, EHR) and multimodal models (MedFuse, DrFuse, and MeTRA) across all 25 clinical conditions. These modality-level grids highlight how calibration dynamics evolve with selective rejection and emphasize that miscalibration, particularly in positive predictions, varies across architectures but follows the same underlying trend.

For the unimodal CXR model, positive-class ECE exhibits substantial variance across seeds, reflecting unstable confidence estimates despite selective AUC patterns that parallel those of MedFuse. The unimodal EHR model shows stratified calibration curves closely aligned with the multimodal results (Figure 2.a), though with consistently higher positive-class ECE for several low-prevalence conditions.

Extending this analysis to DrFuse and MeTRA reveals broadly similar behaviors: both models produce the same selective and calibration trends observed in MedFuse, with only minor condition-specific deviations. These differences, however, do not translate into systematic performance or calibration gains, suggesting that architectural modifications, whether dynamic fusion or transformer-based attention, do not substantially alter the fundamental reliability dynamics across modalities.

**Unimodal CXR**



Figure C5: **Stratified calibration and selective prediction for unimodal CXR.** For unimodal CXR, stratified positive-class ECE exhibits substantial variability across seeds and conditions (top), and AUROC/AUPRC remain relatively stable but clearly below the performance of the multimodal MedFuse model. Across conditions, higher positive-class ECE AUC is associated with lower selective AUC (middle), and we reject the null hypothesis of no non-negative association between stratified ECE and Selective AUC (Spearman's rank $p < 0.05$), mirroring the instability observed in the multimodal setting. In contrast, negative-class ECE AUC shows no clear monotonic relationship with selective AUC (bottom), with points spread broadly along the ECE axis, indicating weaker and less interpretable connections between negative-class calibration and selective prediction.

**Unimodal EHR**



Figure C6: **Stratified calibration and selective prediction for unimodal EHR.** For unimodal EHR, positive-class ECE is higher for several conditions, indicating that EHR-driven predictions are more prone to overconfidence on minority-positive cases. However, both stratified ECE and AUROC/AUPRC exhibit noticeably less variability across seeds and conditions than CXR, suggesting more stable (less noisy) predictions; the resulting trends closely mirror those observed in MedFuse, consistent with EHR driving much of the multimodal performance. In the correlation analysis, positive-class ECE AUC and selective AUC remain negatively associated, though point clusters show more abrupt shifts along the ECE axis and slightly more spread than for CXR. As with CXR, negative-class ECE AUC shows no clear monotonic relationship with Selective AUC, and no consistent trend emerges from the scatter.

**MedFuse**



Figure C7: **Stratified calibration and selective prediction for MedFuse.** Across conditions, higher positive-class ECE AUC consistently predicts lower selective AUC, indicating that overconfidence in positive predictions drives instability in selective evaluation. In contrast, most conditions cluster at low negative-class ECE values with positive selective AUC, and no clear monotonic relationship emerges, indicating that well-calibrated negative predictions reliably support selective gains, while miscalibration is primarily driven by overconfident positive predictions.

**DrFuse**



Figure C8: **Stratified calibration and selective prediction for DrFuse.** For DrFuse, stratified ECE and selective AUROC/AUPRC across 25 conditions closely mirror the patterns observed in MedFuse, with only modest gains in aggregate AUROC and AUPRC. As in the other models, higher positive-class ECE AUC is associated with lower Selective AUC, indicating that overconfident positive predictions continue to drive instability in selective evaluation. Negative-class ECE AUC again shows no clear monotonic relationship with Selective AUC, reinforcing that changing the fusion architecture alone does not substantially improve calibration or selective prediction reliability.

**MeTRA**



Figure C9: **Stratified calibration and selective prediction for MeTRA.** For MeTRA, stratified ECE and selective AUROC/AUPRC across 25 conditions show overall performance comparable to MedFuse and DrFuse, with similar class-dependent calibration issues. However, positive-class ECE exhibits noticeably higher variance across seeds and slightly noisier selective AUROC/AUPRC trends, which may reflect the adaptation of MeTRA from its original in-hospital mortality task to clinical condition classification. As with the other multimodal architectures, negative-class ECE AUC shows no clear monotonic relationship with selective AUC, indicating that architectural complexity alone does not resolve miscalibration-driven instability in selective prediction.

## 5. Loss Upweight Training of Multimodal Architectures

To assess whether simple training interventions can mitigate the reliability issues observed in our study, we apply label-dependent loss upweighting to each multimodal backbone (MedFuse, DrFuse, and MeTRA) and analyze the resulting stratified calibration and selective prediction behavior.

**MedFuse Upscaled**



Figure C10: **Stratified calibration and selective prediction for MedFuse (loss-upweighted).** With label-dependent loss upweighting, MedFuse shows reduced variance and lower positive-class ECE across many conditions, although new calibration spikes appear at higher rejection thresholds and selective AUROC/AUPRC remain non-monotonic, indicating that the model is still not well calibrated. Negative-class ECE also worsens for some conditions despite stable discrimination metrics. In the correlation analysis, positive-class ECE AUC and Selective AUC remain negatively associated, with points now clustered more tightly at lower ECE values; for negative-class ECE, a clearer negative relationship with selective AUC emerges and the null of no non-negative association is rejected for both metrics, despite substantial noise in AUPRC. Overall, loss upweighting pulls conditions toward a lower-ECE regime but does not fully resolve the underlying miscalibration driving selective prediction instability.
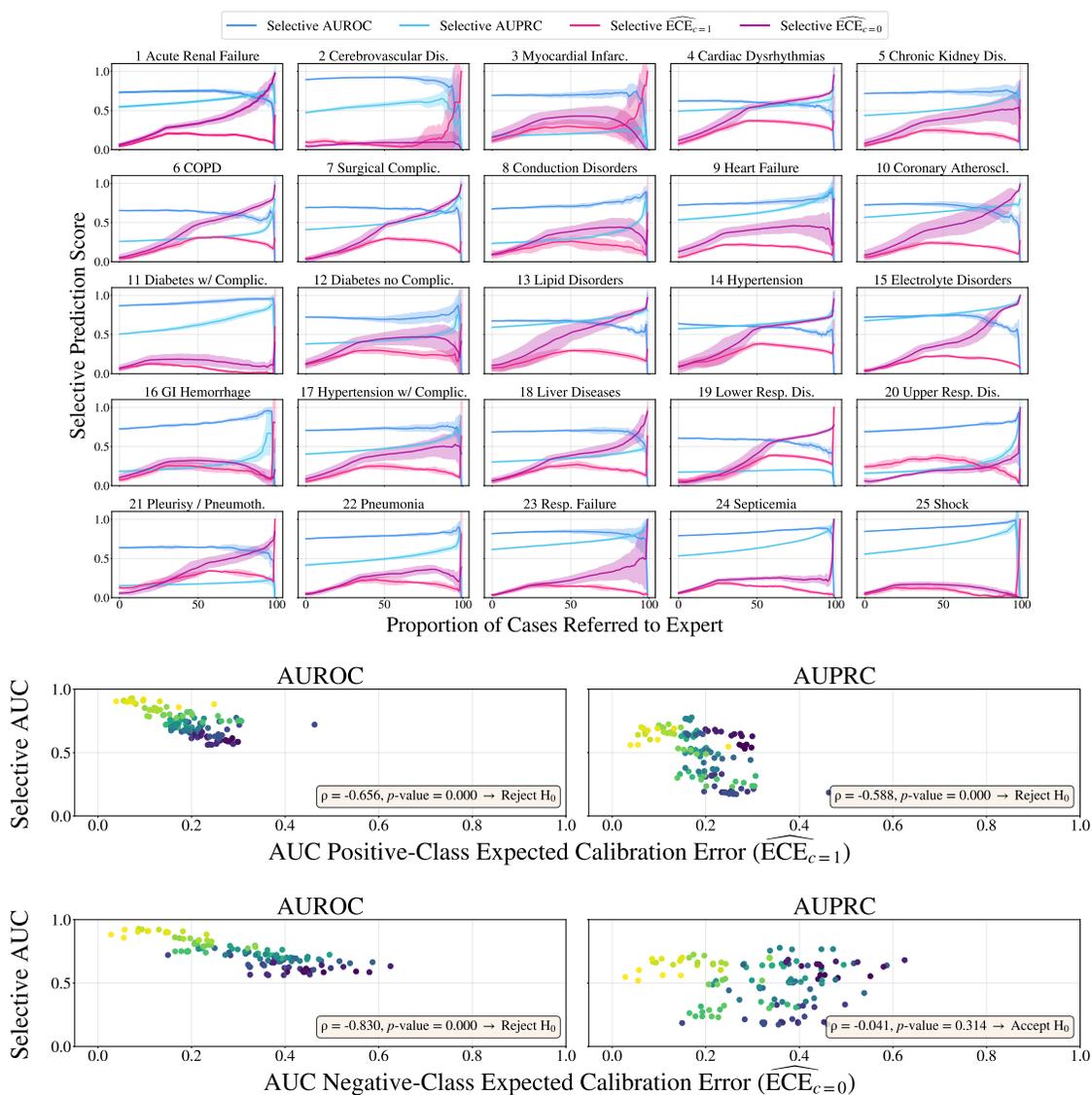
**DrFuse Upscaled**



Figure C11: **Stratified calibration and selective prediction for DrFuse (loss-upweighted).** For loss-upweighted DrFuse, stratified ECE and selective AUROC/AUPRC show trends broadly similar to MedFuse, but with smaller apparent gains and substantially higher variance across seeds for negative-class ECE. Across conditions, positive-class ECE AUC remains negatively associated with Selective AUC, indicating that overconfident positive predictions still drive selective instability. For negative-class ECE, however, the correlation analysis is less conclusive: AUROC exhibits a weak negative trend, but AUPRC is highly noisy and the null hypothesis of no non-negative association is not consistently rejected, suggesting that loss upweighting yields limited and architecture-dependent improvements in calibration.
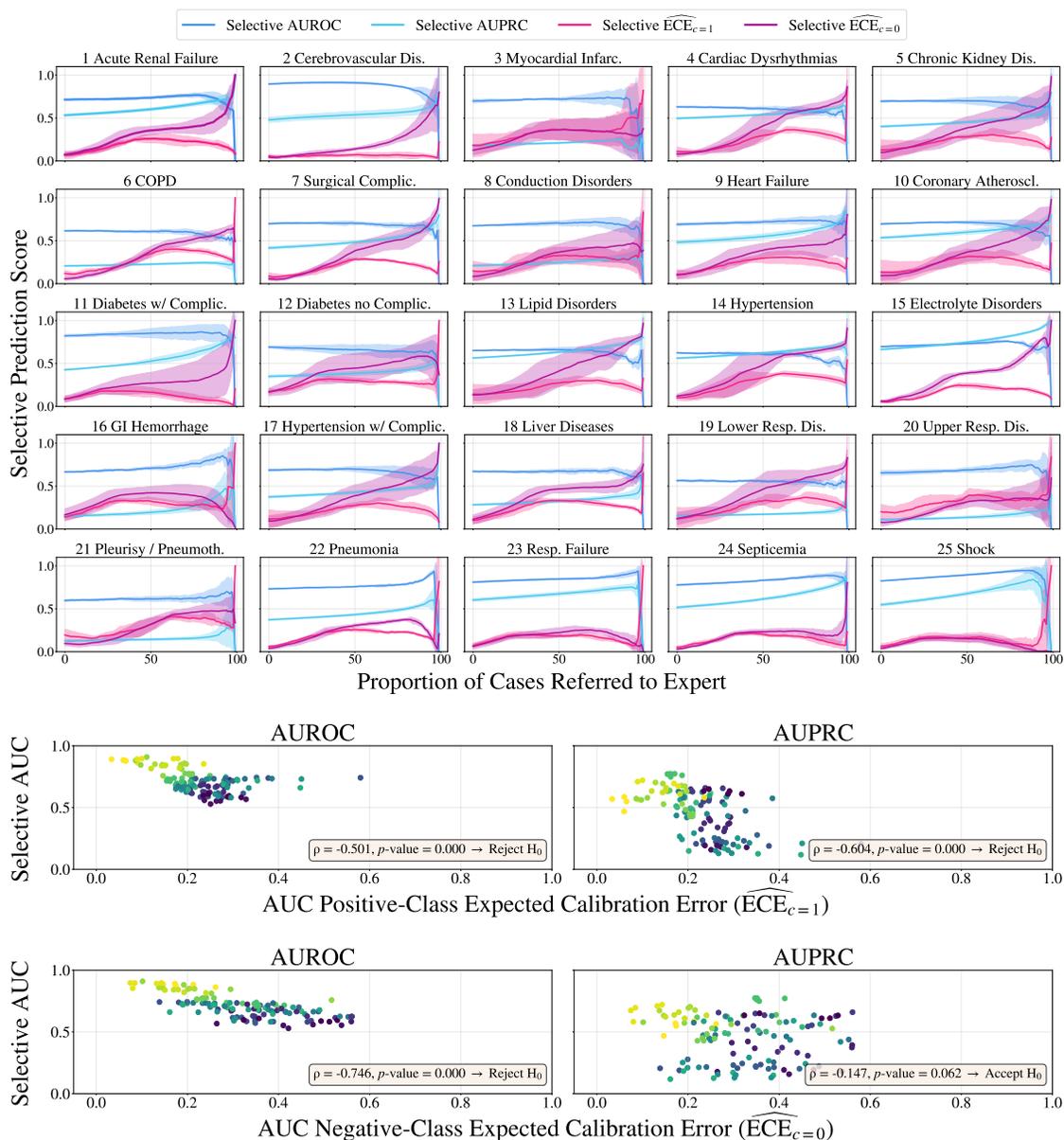
**MeTRA Upscaled**



Figure C12: **Stratified calibration and selective prediction for MeTRA (loss-upweighted).** For loss-upweighted MeTRA, stratified ECE and selective AUROC/AUPRC exhibit improvements similar to those observed in MedFuse and DrFuse, with reduced positive-class ECE for several conditions but persistent variability in negative-class ECE. The correlation analysis again shows a negative association between positive-class ECE AUC and Selective AUC, while negative-class ECE displays noisier and less stable trends, particularly for AUPRC, reinforcing that label-dependent loss upweighting yields only partial and architecture-consistent gains in calibration and selective prediction reliability.

# 6. Per-Condition Comparative Analysis

To provide qualitative insight into how individual modalities and uncertainty-aware training influence selective prediction, we examine three representative conditions: Acute Cerebrovascular Disease (2), Conduction Disorders (8), and Shock (25), across all model variants (EHR, CXR, MedFuse, MedFuse (Loss-Upweight), DrFuse and MeTra). These examples illustrate how modality-specific calibration behaviors interact and how group-aware priors (GAP) mitigate instability and positive-class miscalibration observed in prior sections.

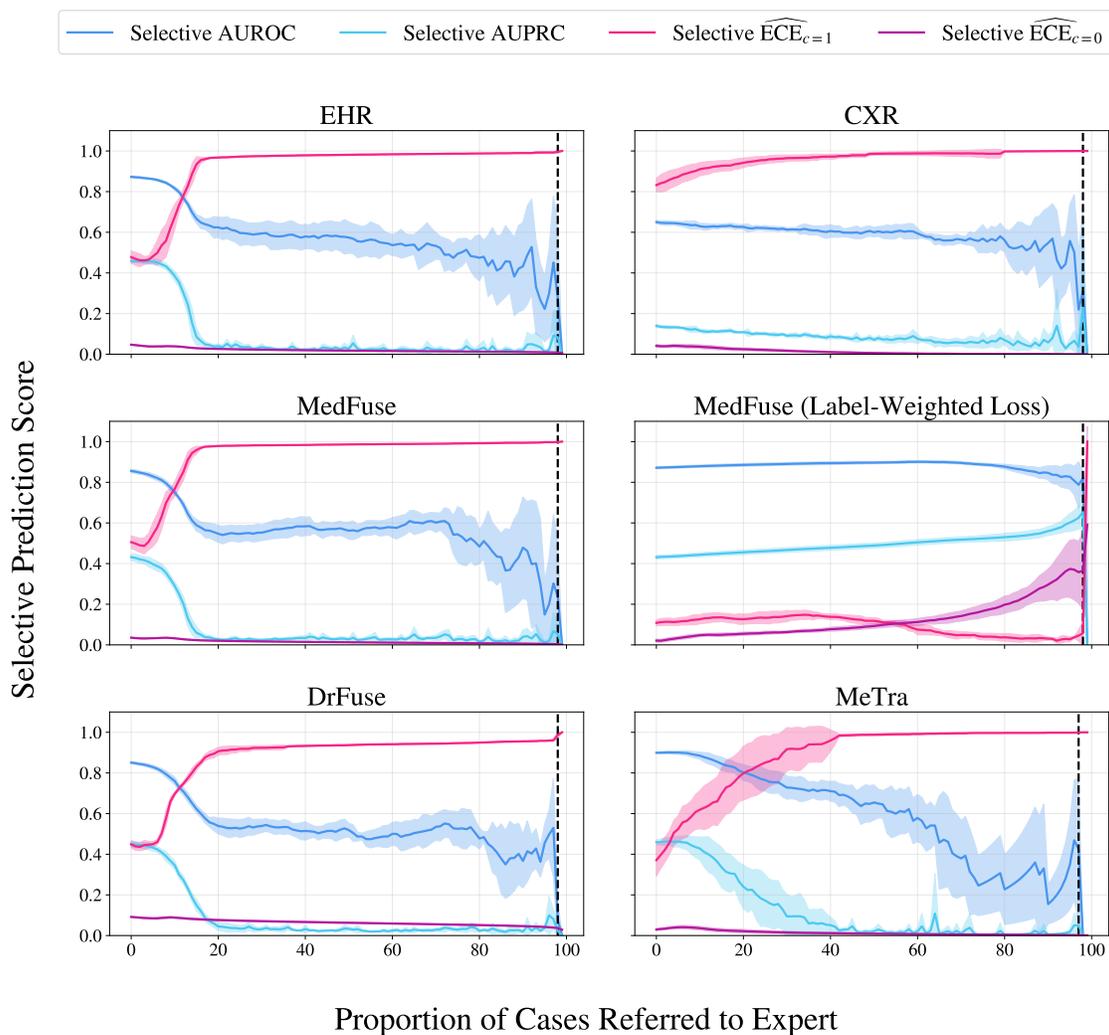### 6.1. Label 2: Acute Cerebrovascular Disease



Figure D1: **Comparative analysis for Acute Cerebrovascular Disease.** For this condition, EHR predominantly drives the discriminative performance of MedFuse, while CXR shows poor positive-class calibration and contributes to the oscillatory behavior of selective AUROC. Applying label-weighted loss reduces early-threshold miscalibration in positive predictions but does not fully correct calibration at higher thresholds, leaving the overall selective prediction trend largely unchanged.
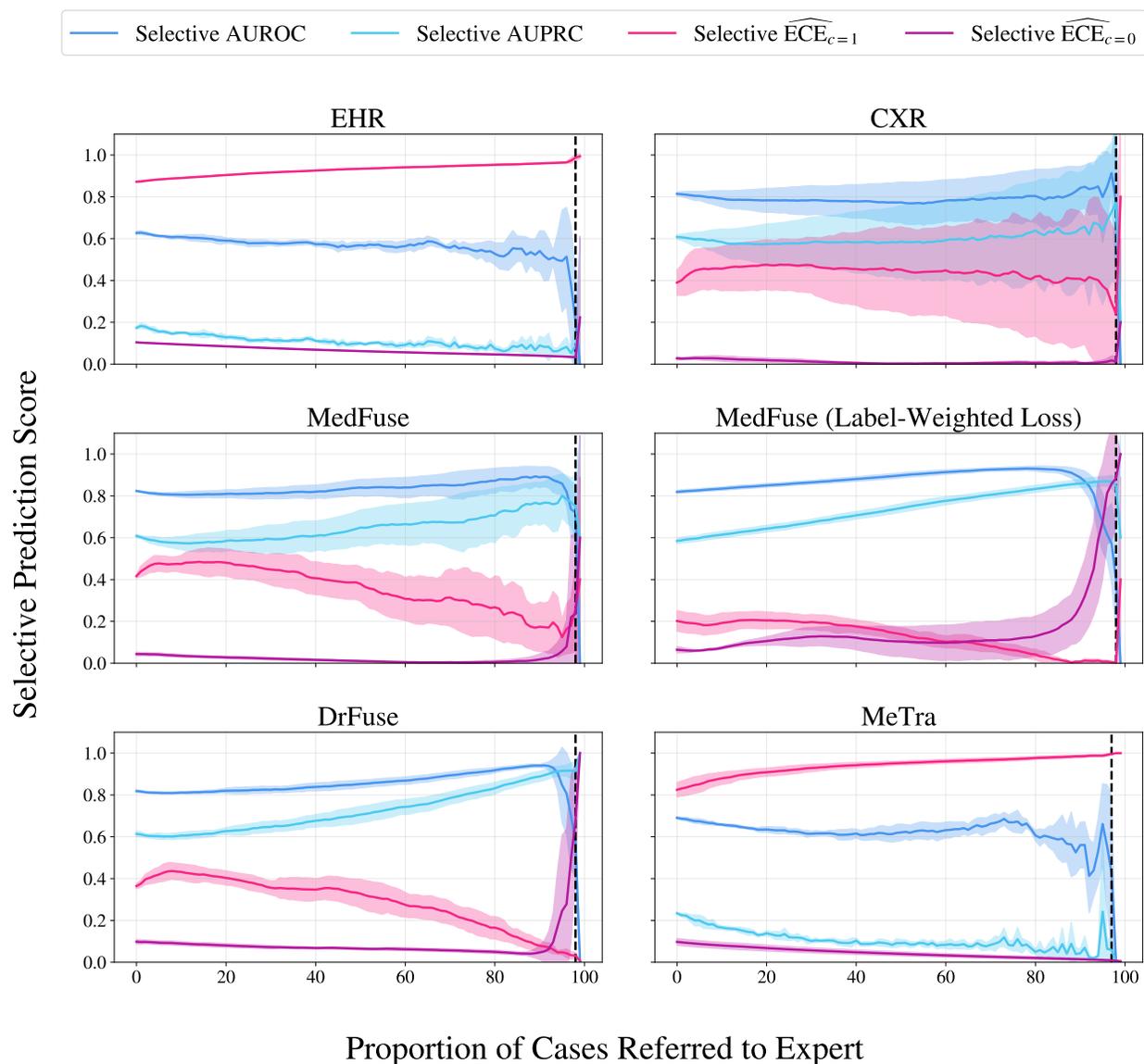
## 6.2. Label 8: Conduction Disorders



Figure D2: **Comparative analysis for Conduction Disorders.** For this condition, the CXR model shows substantial seed-dependent variability, achieving higher mean performance but poorer calibration consistency than EHR. MedFuse integrates both modalities to improve selective AUROC and stabilize performance, though positive-class ECE remains variable. Label-weighted loss further reduces seed variance and attenuates positive-class ECE at higher confidence thresholds, yielding smoother and more reliable selective behavior, but with similar performance drops at the final thresholds.
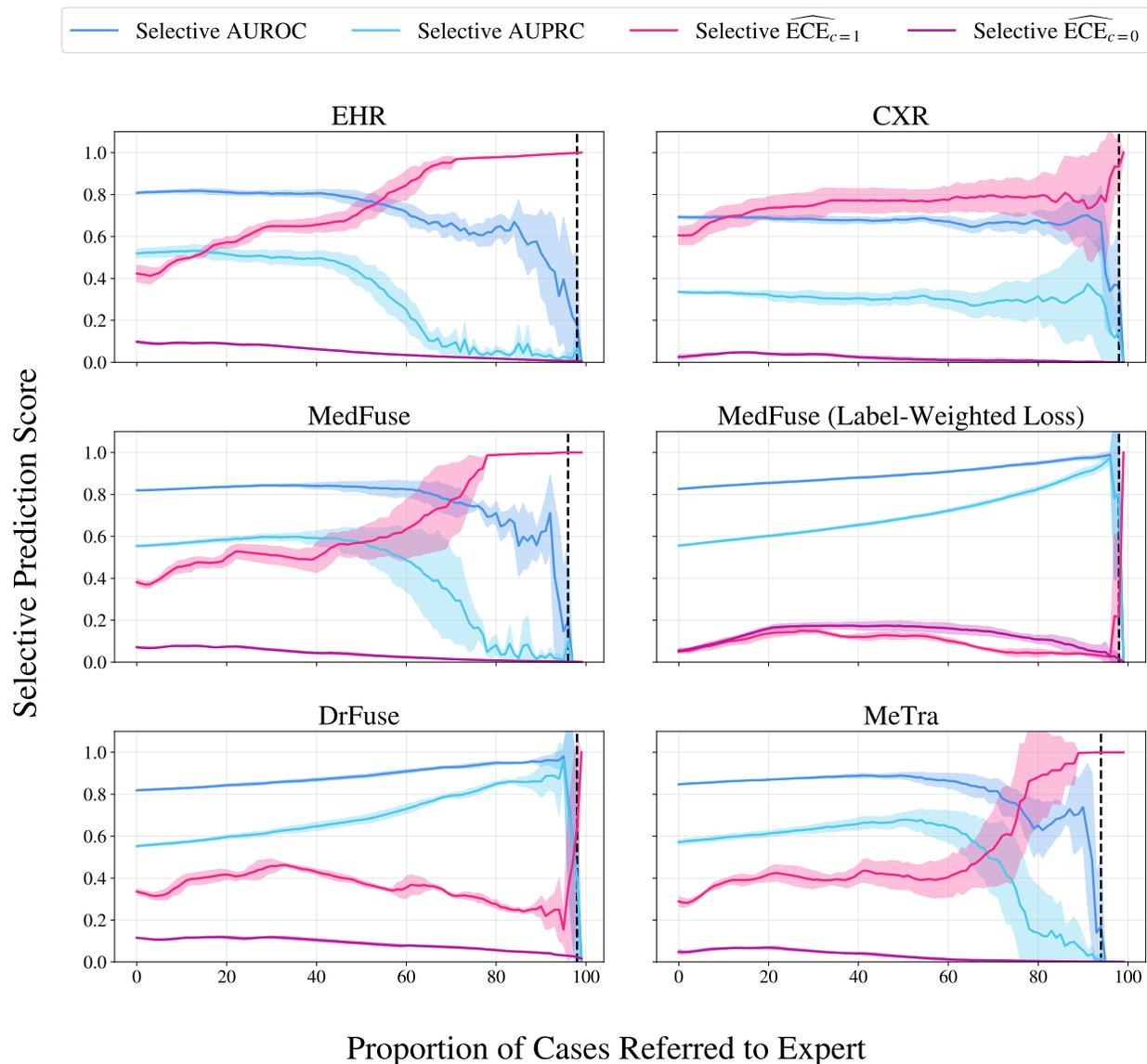
## 6.3. Label 25: Shock



Figure D3: **Comparative analysis for Shock.** For this condition, unimodal EHR and CXR achieve comparable baseline performance, with EHR showing slightly lower positive-class ECE at early thresholds. As rejection increases, however, EHR's AUROC quickly deteriorates while CXR remains uniformly undercalibrated, and MedFuse largely inherits these calibration inconsistencies. Label-weighted loss substantially improves both stability and calibration, suppressing positive-class ECE spikes and producing a clearer monotonic gain in selective AUROC. This condition provides one of the strongest examples of how simple loss upweighting can directly reduce positive-class miscalibration and enhance the reliability of uncertainty estimates.