# Informative Priors Improve the Reliability of Multimodal Clinical Data Classification

**L. Julian Lechuga Lopez**[1,2]
**Tim G. J. Rudner**[2]
**Farah E. Shamout**[1,2]
[1]*NYU Abu Dhabi, AD, UAE*
[2]*New York University, NY, USA*

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

## Abstract

Machine learning-aided clinical decision support has the potential to significantly improve patient care. However, existing efforts in this domain for principled quantification of uncertainty have largely been limited to applications of ad-hoc solutions that do not consistently improve reliability. In this work, we consider stochastic neural networks and design a tailor-made multimodal data-driven (M2D2) prior distribution over network parameters. We use simple and scalable Gaussian mean-field variational inference to train a Bayesian neural network using the M2D2 prior. We train and evaluate the proposed approach using clinical time-series data in MIMIC-IV and corresponding chest X-ray images in MIMIC-CXR for the classification of acute care conditions. Our empirical results show that the proposed method produces a more reliable predictive model compared to deterministic and Bayesian neural network baselines.

**Keywords:** Uncertainty quantification, multimodal healthcare data, Bayesian inference

## 1. Introduction

Trustworthy machine learning in healthcare requires robust uncertainty quantification (Begoli et al., 2019; Gruber et al., 2023), considering the safety-critical nature of clinical practice. Sources of uncertainty can be due to model parameters, noise and bias of the calibration data, or deployment of the model in an out-of-distribution scenario (Miller et al., 2014).

Unfortunately, the literature in machine learning for healthcare has largely neglected developing tailored solutions for improved uncertainty quantification (Kompa et al., 2021), perhaps due to the limited underlying theory on how to best adapt predictive un-

certainty in clinical tasks (Begoli et al., 2019). Other challenges include the complexity of scaling uncertainty quantification in real-time clinical systems, limited empirical evaluation of different methods due to the lack of well-constructed priors by medical experts (Zou et al., 2023), and the high prevalence of data shifts in real-world clinical applications that can negatively affect predictive performance (Ovadia et al., 2019b; Xia et al., 2022), further emphasizing the need for better uncertainty in predictive models.

Additionally, despite the recent proliferation of multimodal learning, existing work on uncertainty quantification in healthcare has mainly been studied in the unimodal setting, with a particular focus on medical imaging applications (Gawlikowski et al., 2021). This includes brain tumor segmentation (Jungo et al., 2018), skin lesion segmentation (DeVries and Taylor, 2018), and diabetic retinopathy detection tasks (Filos et al., 2019; Band et al., 2021; Nado et al., 2022), among others. Hence, effective quantification of predictive uncertainty in the context of multimodal clinical problems remains a challenging and unsolved task (Tran et al., 2022).

We propose a multimodal data-driven (M2D2) prior distribution over neural network parameters to improve uncertainty quantification in multimodal fusion of chest X-ray images and clinical time series data. We evaluate the use of effective priors on the two unimodal components of the multimodal fusion network: an image-based convolutional neural network and a recurrent neural network for clinical time series. In summary, we make the following contributions:

1. We design a multimodal data-driven (M2D2) prior distribution over neural network parameters that places high probability density on desired predictive functions.
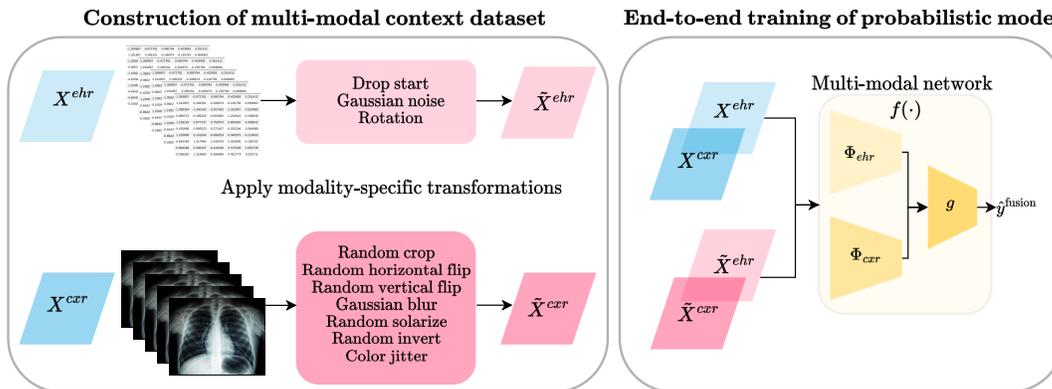
Figure 1: **Overview of model training.** *Left:* We construct a multimodal context dataset by applying modality-specific transformations to clinical time series data and chest X-ray images, resulting in a tailored distribution shift. *Right:* We train the multimodal neural network end-to-end with the training set and the constructed context dataset.

2. We evaluate the method on large publicly-available multimodal datasets: MIMIC-IV and MIMIC-CXR (Johnson et al., 2019, 2021), for the classification of acute care conditions assigned to patient stays in the intensive care unit.

3. Our findings illustrate an increase in predictive performance and improved reliability in uncertainty-aware selective prediction.

## 2. Related Work

### 2.1. Multimodal Learning in Healthcare

Multimodal learning in healthcare seeks to exploit complementary information from different data modalities to enhance the predictive capabilities of learning models. There are different approaches for leveraging information across different data modalities, with the most popular paradigm being multimodal fusion (Huang et al., 2020). For example, Zhang et al. (2020) and Calhoun and Sui (2016) investigated different methods for fusion segmentation and quantification in neuroimaging by leveraging different imaging modalities in the same data pipeline. Another recent study focused on the development of smart healthcare applications by merging multimodal signals collected from different types of medical sensors (Muhammad et al., 2021). Other studies also show improved predictive performance when using multiple modalities in prognostic tasks in patients with COVID-19 (Shamout et al., 2021; Jiao et al., 2021).

Despite the promise of multimodal learning in healthcare, research in reliable uncertainty quantification applications in the multimodal setting is currently limited. There is no generalized use of uncertainty quantification methods that address increased data distribution shifts and deal with multiple modalities simultaneously (Liang et al., 2023).

### 2.2. Variational Inference in Neural Networks

We consider a stochastic neural network $f(\cdot\,;\Theta)$, defined in terms of stochastic parameters $\Theta \in \mathbb{R}^P$. For an observation model $p_{Y|X,\Theta}$ and a prior distribution over parameters $p_\Theta$, Bayesian inference provides a mathematical formalism for finding the posterior distribution over parameters given the observed data, $p_{\Theta|\mathcal{D}}$ (MacKay, 1992; Neal, 1996). However, since neural networks are non-linear in their parameters, exact inference over the stochastic network parameters is analytically intractable.

Variational inference is an approach that seeks to avoid this intractability by framing posterior inference as finding an approximation $q_\Theta$ to the posterior $p_{\Theta|\mathcal{D}}$ via the variational optimization problem:

$$\min_{q_\Theta \in \mathcal{Q}_\Theta} D_{\mathrm{KL}}(q_\Theta \parallel p_{\Theta|\mathcal{D}}) \iff \max_{q_\Theta \in \mathcal{Q}_\Theta} \mathcal{F}(q_\Theta),$$

where $\mathcal{F}(q_\Theta)$ is the variational objective

$$\mathcal{F}(q_\Theta) \doteq \mathbb{E}_{q_\Theta}[\log p(y_\mathcal{D} \mid x_\mathcal{D}, \Theta)] - D_{\mathrm{KL}}(q_\Theta \parallel p_\Theta), \quad (1)$$

$\mathcal{Q}_\Theta$ is a variational family of distributions (Wainwright and Jordan, 2008), and $(x_\mathcal{D}, y_\mathcal{D})$ are the training data. One particularly simple type of variational inference is Gaussian mean-field variational inference (Blundell et al., 2015; Graves, 2011), where the posterior distribution over network parameters is approximated by a Gaussian distribution with a diagonal covariance matrix. This method enables stochastic optimization and
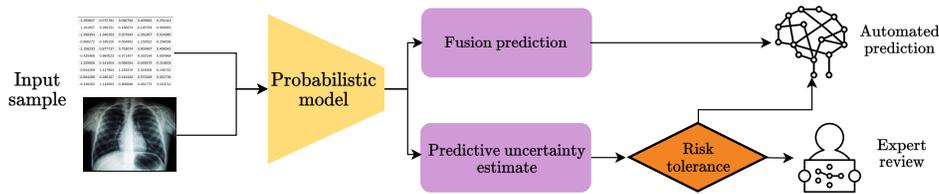
Figure 2: **Selective prediction algorithm.** For each input sample, we obtain a prediction and an estimate of the model's uncertainty for that specific data point. If the uncertainty estimate is higher than the selected risk tolerance, then the sample will be sent to an expert for further review and classification. Otherwise, the sample is processed by the learning model for automated prediction.

can be scaled to large neural networks (Hoffman et al., 2013). However, Gaussian mean-field variational inference has been shown to underperform with deterministic neural networks when uninformative, standard Gaussian priors are used (Ovadia et al., 2019a; Rudner et al., 2022a).

To improve the performance, we extend the approach presented in Rudner et al. (2023a) to stochastic neural networks, construct a data-driven prior distribution from multimodal input data, and use this prior for Gaussian mean-field variational inference to improve the performance of neural networks for multimodal clinical prediction tasks.

## 3. Constructing Data-Driven Priors for Models with Multimodal Input Data

We consider a supervised multimodal fusion task on data $\mathcal{D} \doteq \{(x_n^1, x_n^2, y_n^{\text{fusion}})\}_{n=1}^N = (X_{\mathcal{D}}^1, X_{\mathcal{D}}^2, Y_{\mathcal{D}})$. As shown in Figure 1, we consider the first modality to be clinical time series data extracted from electronic health records, denoted by $X^{\text{ehr}}$, and the second to be chest X-ray images, denoted by $X^{\text{cxr}}$. For a given sample $(x^{\text{ehr}}, x^{\text{cxr}})$, the two modalities are processed by encoders $\Phi_{\text{ehr}}$ and $\Phi_{\text{cxr}}$ respectively, and their concatenated feature representations are further processed by a classifier $g(\cdot)$ and activation function to compute the fusion prediction $\hat{y}^{\text{fusion}}$. The loss is computed based on the predictions and ground truth labels $y^{\text{fusion}} \in \mathcal{Y}$, where $\mathcal{Y} \subseteq \{0, 1\}^Q$, with $Q > 1$ for multi-label classification.

### 3.1. Informative Priors for Multimodal Data

One of the key components in defining the probabilistic model for our uncertainty quantification method is the definition of a sensible and explainable prior distribution. In this work, we construct a prior distribution over parameters that places high probability density on parameter values that induce predictive functions that have high uncertainty on input points that are meaningfully different from the training data. To do this, we build on the approach proposed in Rudner et al. (2023a) and use information about the two input modalities to construct a data-driven prior that can help find an approximate posterior distribution with desirable properties (e.g., an induced predictive distribution with reliable uncertainty estimation). More specifically, we construct a data-driven prior over some set of model parameters $\Psi$ and condition it on a set of context points $\tilde{X}$, that is $p(\psi|\tilde{x})$. In Appendix A, we show that we can derive a tractable variational objective using this prior. The objective is given by Equation (A.16).

To construct a meaningful prior, we need to specify a distribution over the set of context points, $p_{\tilde{X}}$. We design a multimodal prior by letting $\tilde{X}$ be a set of randomly generated multimodal input points $(\tilde{X}^{\text{ehr}}, \tilde{X}^{\text{cxr}})$ designed to be distinct from the training data. For the clinical time series data, we construct $\tilde{X}^{\text{ehr}}$ by applying three transformations to the original time series: drop start, Gaussian noise, and inversion (i.e, for each $x_i$ in $1, ..., n$, $x_1 = x_n$, $x_2 = x_{n-1}$, $x_3 = x_{n-2}$, etc). For the chest X-ray images, we construct $\tilde{X}^{\text{cxr}}$ by applying seven transformations representative of perturbations that exist in real-world medical settings to the imaging data: random crop, random horizontal and vertical flip, Gaussian blur, random solarize, random invert and color jitter.

Hence, this context set encompasses distributionally shifted points, where we want the model's uncertainty to be higher.

## 4. Empirical Evaluation

To evaluate the proposed approach, we combine clinical time series data from MIMIC-IV (Johnson et al., 2021) and chest X-ray images from MIMIC-CXR

Table 1: **Performance results.** We summarize the results on the test set for the baselines and our stochastic model, including 95% confidence intervals computed via bootstrapping. Higher values are better for all metrics.

| Model (MedFuse) | AUROC | AUPRC | Selective AUROC | Selective AUPRC |
|---|---|---|---|---|
| **Deterministic** (Hayat et al., 2022) | 0.726 (0.718, 0.733) | 0.503 (0.493, 0.517) | 0.724 (0.715, 0.735) | 0.439 (0.429, 0.455) |
| **Bayesian** (standard prior) | 0.729 (0.722, 0.736) | 0.507 (0.497, 0.521) | **0.748** (0.739, 0.758) | 0.448 (0.437, 0.467) |
| **Bayesian** (M2D2 prior) (**Ours**) | **0.735** (0.728, 0.742) | **0.514** (0.504, 0.528) | **0.748** (0.738, 0.760) | **0.452** (0.441, 0.472) |

(Johnson et al., 2019) collected during the same patient stay in the intensive care unit for multi-label classification of acute care conditions.

## 4.1. Experimental Setup

We follow the pre-processing steps and use the same neural network architecture (MedFuse) as Hayat et al. (2022). $\Phi_{\text{ehr}}$ is a two-layer LSTM network (Hochreiter and Schmidhuber, 1997), $\Phi_{\text{cxr}}$ is a ResNet-34 (He et al., 2015), $g(\cdot)$ is a fully connected layer, and $\hat{y}^{\text{fusion}}$ are the class probabilities obtained by applying a sigmoid function to $g$. We use the paired dataset, such that each sample contains both modalities (i.e., there are no missing modalities). Hence, the training, validation, and test sets consisted of 7756, 877, and 2161 samples, respectively. We construct the context dataset using the training set.

We train the multimodal network for 400 epochs using the loss presented in Equation (B.17), with the Adam optimizer, a batch size of 16, and a learning rate of $2 \times 10^{-4}$. Further details on the experimental setup and grid-based hyperparameter tuning can be found in Appendix B.

## 4.2. Evaluation Metrics

We evaluate the overall performance of the models on the test set using the AUROC and Area Under the Precision-Recall curve (AUPRC) (Hayat et al., 2022).

In addition, we compute selective prediction evaluation metrics to better assess models' predictive uncertainty. As shown in Figure 2, selective prediction modifies the standard prediction pipeline by introducing a "reject option", $\bot$, via a gating mechanism defined by selection function $s : \mathcal{X} \to \mathbb{R}$ that determines whether a prediction should be made for a given input point $x \in \mathcal{X}$ (El-Yaniv et al., 2010). For rejection threshold $\tau$, with $s$ representing the entropy of $x$, the prediction model is given by

$$(p(y \,|\, \cdot, \theta; f), s)(x) = \begin{cases} p(y \,|\, x, \theta; f), & \text{if } s \leq \tau \\ \bot, & \text{otherwise} \end{cases} \quad (2)$$

To evaluate the predictive performance of a prediction model $(p(y \,|\, \cdot, \theta; f), s)(x)$ with a single label,

we compute the AUROC and AUPRC over rejection thresholds $\tau = 0\%, ..., 99\%$. We then average the metrics across all thresholds, yielding selective prediction AUROC and AUPRC scores that explicitly incorporate both a model's predictive performance and its predictive uncertainty. For our multi-label classification task, we report the average selective prediction scores across all 25 labels.

## 4.3. Results

Table 1 summarizes the performance results on the test set. Additional per-label results are shown in Appendix C. The Bayesian neural network with an M2D2 prior achieves a better AUROC and AUPRC of 0.735 and 0.514, respectively, compared with 0.726 and 0.503 by the deterministic model. It also achieves a higher selective AUROC and AUPRC of 0.748 and 0.452, respectively, compared with 0.724 and 0.439 by the deterministic model. Our proposed approach achieves a comparable selective AUROC to when using a standard prior.

We also observe a decrease in selective AUPRC compared to the 0%-rejection AUPRC. This can occur when a model is poorly calibrated: When the AUPRC for any rejection threshold is below the 0%-rejection score, the selective AUPRC can be lower than the 0%-rejection score. Overall, the selective prediction scores reflect the model's ability to identify samples that are more likely to be misclassified and should be reviewed by a clinician, and as such, are valuable in assessing model reliability in clinical settings.

## 5. Conclusion

We designed a multimodal data-driven (M2D2) prior to improve the reliability of multimodal fusion of clinical time series data and chest X-ray images. We demonstrated that Bayesian neural networks with such a prior achieve better performance, in terms of AUROC, AUPRC, and selective prediction scores, than deterministic models. For future work, we aim to evaluate the proposed approach in settings of missing modalities, on additional tasks, such as in-hospital mortality prediction, and other multi-modal datasets.

## 6. Acknowledgements

## References

Neil Band, Tim G. J. Rudner, Qixuan Feng, Angelos Filos, Zachary Nado, Michael W. Dusenberry, Ghassen Jerfel, Dustin Tran, and Yarin Gal. Benchmarking Bayesian deep learning on diabetic retinopathy detection tasks. 2021.

Edmon Begoli, Tanmoy Bhattacharya, and Dimitri Kusnezov. The need for uncertainty quantification in machine-assisted medical decision making. *Nature Machine Intelligence*, 1(1):20–23, 2019.

Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. volume 37 of *Proceedings of Machine Learning Research*, pages 1613–1622, Lille, France, 07–09 Jul 2015. PMLR.

James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018.

Vince D Calhoun and Jing Sui. Multimodal fusion of brain imaging data: a key to finding the missing link (s) in complex mental illness. *Biological psychiatry: cognitive neuroscience and neuroimaging*, 1(3):230–244, 2016.

Terrance DeVries and Graham W Taylor. Leveraging uncertainty estimates for predicting segmentation quality. *arXiv preprint arXiv:1807.00502*, 2018.

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

Angelos Filos, Sebastian Farquhar, Aidan N Gomez, Tim GJ Rudner, Zachary Kenton, Lewis Smith, Milad Alizadeh, Arnoud De Kroon, and Yarin Gal. A systematic comparison of Bayesian deep learning robustness in diabetic retinopathy tasks. *arXiv preprint arXiv:1912.10481*, 2019.

Jakob Gawlikowski, Cedrique Rovile Njieutcheu Tassi, Mohsin Ali, Jongseok Lee, Matthias Humt, Jianxiang Feng, Anna Kruspe, Rudolph Triebel, Peter Jung, Ribana Roscher, et al. A survey of uncertainty in deep neural networks. *arXiv preprint arXiv:2107.03342*, 2021.

Irving John Good. Rational decisions. *Journal of the Royal Statistical Society: Series B (Methodological)*, 14(1):107–114, 1952.

Alex Graves. Practical variational inference for neural networks. *Advances in neural information processing systems*, 24, 2011.

Cornelia Gruber, Patrick Oliver Schenk, Malte Schierholz, Frauke Kreuter, and Göran Kauermann. Sources of uncertainty in machine learning – a statisticians' view, 2023.

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, 05–06 Aug 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.

Matthew D. Hoffman, David M. Blei, Chong Wang, and John Paisley. Stochastic variational inference. *Journal of Machine Learning Research*, 14(1):1303–1347, May 2013. ISSN 1532-4435.

Shih-Cheng Huang, Anuj Pareek, Saeed Seyyedi, Imon Banerjee, and Matthew P Lungren. Fusion of medical imaging and electronic health records using deep learning: a systematic review and implementation guidelines. *NPJ digital medicine*, 3(1):136, 2020.

Zhicheng Jiao, Ji Whae Choi, Kasey Halsey, Thi My Linh Tran, Ben Hsieh, Dongcui Wang, Feyisope Eweje, Robin Wang, Ken Chang, Jing Wu, et al. Prognostication of patients with covid-19 using artificial intelligence based on chest x-rays

and clinical data: a retrospective study. *The Lancet Digital Health*, 3(5):e286–e294, 2021.

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.

Alain Jungo, Richard McKinley, Raphael Meier, Urspeter Knecht, Luis Vera, Julián Pérez-Beteta, David Molina-García, Víctor M Pérez-García, Roland Wiest, and Mauricio Reyes. Towards uncertainty-assisted brain tumor segmentation and survival prediction. In *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries: Third International Workshop, BrainLes 2017, Held in Conjunction with MICCAI 2017, Quebec City, QC, Canada, September 14, 2017, Revised Selected Papers 3*, pages 474–485. Springer, 2018.

Leo Klarner, Tim G. J. Rudner, Michael Reutlinger, Torsten Schindler, Garrett M. Morris, Charlotte Deane, and Yee Whye Teh. Drug discovery under covariate shift with domain-informed prior distributions over functions. In *Proceedings of the 40th International Conference on Machine Learning*, 2023.

Benjamin Kompa, Jasper Snoek, and Andrew L Beam. Second opinion needed: communicating uncertainty in medical machine learning. *NPJ Digital Medicine*, 4(1):4, 2021.

Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. Foundations and trends in multimodal machine learning: Principles, challenges, and open questions, 2023.

David J. C. MacKay. A practical Bayesian framework for backpropagation networks. *Neural Comput.*, 4(3):448–472, May 1992. ISSN 0899-7667. doi: 10.1162/neco.1992.4.3.448.

David C Miller, Brenda Ng, John Eslick, Charles Tong, and Yang Chen. Advanced computational tools for optimization and uncertainty quantification of carbon capture processes. In *Computer Aided Chemical Engineering*, volume 34, pages 202–211. Elsevier, 2014.

Ghulam Muhammad, Fatima Alshehri, Fakhri Karray, Abdulmotaleb El Saddik, Mansour Alsulaiman, and Tiago H Falk. A comprehensive survey on multimodal medical signals fusion for smart healthcare systems. *Information Fusion*, 76:355–375, 2021.

Zachary Nado, Neil Band, Mark Collier, Josip Djolonga, Michael W. Dusenberry, Sebastian Farquhar, Qixuan Feng, Angelos Filos, Marton Havasi, Rodolphe Jenatton, Ghassen Jerfel, Jeremiah Liu, Zelda Mariet, Jeremy Nixon, Shreyas Padhy, Jie Ren, Tim G. J. Rudner, Faris Sbahi, Yeming Wen, Florian Wenzel, Kevin Murphy, D. Sculley, Balaji Lakshminarayanan, Jasper Snoek, Yarin Gal, and Dustin Tran. Uncertainty baselines: Benchmarks for uncertainty & robustness in deep learning, 2022.

Radford M Neal. Bayesian Learning for Neural Networks. 1996.

Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D. Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model's uncertainty? Evaluating predictive uncertainty under dataset shift. In *Advances in Neural Information Processing Systems 32*. 2019a.

Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model's uncertainty? evaluating predictive uncertainty under dataset shift. *Advances in neural information processing systems*, 32, 2019b.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.

Tim G. J. Rudner, Zonghao Chen, Yee Whye Teh, and Yarin Gal. Tractable function-space variational inference in Bayesian neural networks. In *Advances in Neural Information Processing Systems 35*, 2022a.

Tim G. J. Rudner, Freddie Bickford Smith, Qixuan Feng, Yee Whye Teh, and Yarin Gal. Continual learning via sequential function-space variational inference. In *Proceedings of the 39th International Conference on Machine Learning*, 2022b.

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.

Farah E Shamout, Yiqiu Shen, Nan Wu, Aakash Kaku, Jungkyu Park, Taro Makino, Stanislaw Jastrzkeski, Jan Witowski, Duo Wang, Ben Zhang, et al. An artificial intelligence system for predicting the deterioration of covid-19 patients in the emergency department. *NPJ digital medicine*, 4(1):80, 2021.

Dustin Tran, Jeremiah Liu, Michael W. Dusenberry, Du Phan, Mark Collier, Jie Ren, Kehang Han, Zi Wang, Zelda Mariet, Huiyi Hu, Neil Band, Tim G. J. Rudner, Karan Singhal, Zachary Nado, Joost van Amersfoort, Andreas Kirsch, Rodolphe Jenatton, Nithum Thain, Honglin Yuan, Kelly Buchanan, Kevin Murphy, D. Sculley, Yarin Gal, Zoubin Ghahramani, Jasper Snoek, and Balaji Lakshminarayanan. Plex: Towards reliability using pretrained large model extensions, 2022.

Martin J Wainwright and Michael I Jordan. *Graphical Models, Exponential Families, and Variational Inference*. Now Publishers Inc., Hanover, MA, USA, 2008. ISBN 1601981848.

Tong Xia, Jing Han, and Cecilia Mascolo. Benchmarking uncertainty quantification on biosignal classification tasks under dataset shift. In *Multimodal AI in healthcare: A paradigm shift in health intelligence*, pages 347–359. Springer, 2022.

Yu-Dong Zhang, Zhengchao Dong, Shui-Hua Wang, Xiang Yu, Xujing Yao, Qinghua Zhou, Hua Hu, Min Li, Carmen Jiménez-Mesa, Javier Ramirez, et al. Advances in multimodal data fusion in neuroimaging: overview, challenges, and novel orientation. *Information Fusion*, 64:149–187, 2020.

Ke Zou, Zhihao Chen, Xuedong Yuan, Xiaojing Shen, Meng Wang, and Huazhu Fu. A review of uncertainty estimation and its application in medical imaging. *arXiv preprint arXiv:2302.08119*, 2023.

# Supplementary Material

## Appendix A. Variational Objective

Let the mapping $f$ in the parametric observation model $p_{Y|X,\Theta}(y\,|\,x,\theta;f)$ be defined by $f(\cdot\,;\theta)\doteq h(\cdot\,;\theta_h)\theta_L$. For a neural network model, $h(\cdot\,;\theta_h)$ is the post-activation output of the penultimate layer, $\Theta_L$ is the set of stochastic final-layer parameters, $\Theta_h$ is the set of stochastic non-final-layer parameters, and $\Theta\doteq\{\Theta_h,\Theta_L\}$ is the full set of stochastic parameters.

To derive an uncertainty-aware prior distribution over the set of random parameters $\Theta$, we start by specifying an auxiliary inference problem. Let $\tilde{x}=\{x_1,...,x_M\}$ be a set of context points with corresponding labels $\tilde{y}$, and define a corresponding likelihood function $\tilde{p}_{Y|X,\Theta}(\tilde{y}\,|\,\tilde{x},\theta)$ and a prior over the model parameters, $p_\Theta(\theta)$. For notational simplicity, we will drop the subscripts as we advance except when needed for clarity. By Bayes' Theorem, we can write the posterior under the context points and labels as

$$\tilde{p}(\theta\,|\,\tilde{x},\tilde{y})\propto\tilde{p}(\tilde{y}\,|\,\tilde{x},\theta_h)p(\theta_h)p(\theta_L). \tag{A.1}$$

To define a likelihood function that induces a posterior with desirable properties, we start from the same step as Rudner et al. (2023a) and consider the following stochastic linear model for an arbitrary set of points $x\doteq\{x_1,...,x_{M'}\}$,

$$\tilde{Y}_k(x)\doteq h(x;\theta_h)\Theta_k+\varepsilon \quad\text{with}\quad \Theta_k\sim\mathcal{N}(\theta_L;m_k,\tau_f^{-1}s_k) \quad\text{and}\quad \varepsilon\sim\mathcal{N}(\mathbf{0},\tau_f^{-1}\beta I)$$

for output dimensions $k=1,...,K$, where $h(\cdot\,;\theta_h)$ is the feature mapping used to define $f$ evaluated at a set of fixed feature parameters $\theta_h$, $\tau_f$ and $\beta$ are variance parameters, and $m\in\mathbb{R}^{P_L}$ and $s\in\mathbb{R}^{P_L}$ are—for now—fixed parameters for a $P_L$-dimensional final layer. This stochastic linear model induces a distribution over functions (Rudner et al., 2022a,b; Klarner et al., 2023; Rudner et al., 2023b), which—when evaluated at $\tilde{x}$—is given by

$$\mathcal{N}(\tilde{y}_k(\tilde{x});h(\tilde{x};\theta_h)m_k,\tau_f^{-1}K(\tilde{x},\tilde{x};\theta_h,s)_k), \tag{A.2}$$

where

$$K(\tilde{x},\tilde{x};\theta_h,s)_k\doteq h(\tilde{x};\theta_h)(s_kI)h(\tilde{x};\theta_h)^\top+\beta I \tag{A.3}$$

is an $M$-by-$M$ covariance matrix. Viewing this probability density over function evaluations as a likelihood function parameterized by $\theta$, we diverge from Rudner et al. (2023a) and define

$$\tilde{p}(\tilde{y}_k\,|\,\tilde{x},\theta_h)\doteq\mathcal{N}(\tilde{y}_k;h(\tilde{x};\theta_h)m_k,\tau_f^{-1}K(\tilde{x},\tilde{x};\theta_h,s)_k), \tag{A.4}$$

where—unlike in Rudner et al. (2023a)—we do not assume that $m=\mathbf{0}$ and $s=I$. If we define the auxiliary label distribution as $p_{\tilde{Y}\,|\,\tilde{X}}(\tilde{y}\,|\,\tilde{x})\doteq\delta(\{\mathbf{0},...,\mathbf{0}\}-\tilde{y})$, the likelihood $\tilde{p}(\tilde{y}_k\,|\,\tilde{x},\theta_h)$ favors learnable parameters $\theta_h$ for which the induced distribution over functions has a high likelihood of predicting $\mathbf{0}$. Letting

$$\tilde{p}(\tilde{y}\,|\,\tilde{x},\theta)\doteq\prod_{k=1}^K\tilde{p}(\tilde{y}_k\,|\,\tilde{x},\theta,m_k,s_k),$$

and taking the log of the analytically tractable density $\tilde{p}(\tilde{y}\,|\,\tilde{x},\theta;f)$, we obtain

$$\log\tilde{p}(\tilde{y}\,|\,\tilde{x},\theta_h)\propto-\sum_{k=1}^K\frac{\tau_f}{2}(h(\tilde{x};\theta_h)m_k)^\top K(\tilde{x},\tilde{x};\theta_h,s)_k^{-1}h(\tilde{x};\theta_h)m_k,$$

with proportionality up to an additive constant independent of $\theta$. We define

$$\mathcal{J}(\theta, m, s, \tilde{x}, \tilde{y}) \doteq -\sum_{k=1}^{K} \frac{\tau_f}{2} d_M^2(h(\tilde{x}; \theta_h)m_k - \tilde{y}, K(\tilde{x}, \tilde{x}; \theta_h, s)_k) \tag{A.5}$$

where $d_M^2(\Delta, K) \doteq \Delta^\top K^{-1} \Delta$ is the squared Mahalanobis distance for $\Delta = v - w$. We therefore obtain

$$\arg\max_\theta \tilde{p}(\theta \,|\, \tilde{x}, \tilde{y}) = \arg\max_\theta \mathcal{J}(\theta, m, s, \tilde{x}, \tilde{y}) + \log p(\theta)$$

and hence, maximizing $\mathcal{J}(\theta, m, s, \tilde{x}, \tilde{y}) + \log p(\theta)$ with respect to $\theta$ is mathematically equivalent to maximizing the posterior $\tilde{p}(\theta \,|\, \tilde{x}, \tilde{y})$ and leads to functions that are likely under the distribution over functions induced by the neural network mapping while being consistent with the prior over the network parameters.

However, since the parameters $m$ and $s$ are fixed and appear in the auxiliary likelihood function but not in the predictive function $f(\cdot; \theta)$, the objective above is not a good choice if the goal is to find parameters $\theta$ that induce functions that have high predictive uncertainty on the set of context points. To address this shortcoming, we include these parameters in the observation model as the mean and variance parameters of the final-layer parameters in $f(\cdot; \theta)$, treat them as random variables $M$ and $S$, place a prior over them, and ultimately infer an approximate posterior distribution for both.

In particular, we define a prior over the final-layer parameters $\Theta_L$ as

$$p_{\Theta_L}(\theta_L \,|\, m, s) = \mathcal{N}(\theta_L; m, sI) \tag{A.6}$$

and corresponding hyperpriors

$$p_M(m) = \mathcal{N}(m; \mu_0, \tau_0^{-1}I) \tag{A.7}$$

$$p_S(s) = \text{Lognormal}(s; \mathbf{0}, 2\tau_s^{-1}I). \tag{A.8}$$

As before, we will drop subscripts for brevity unless needed for clarity. The full probabilistic model then becomes

$$p(y \,|\, x, \theta_h, \theta_L; f)\, p(\theta, m, s \,|\, \tilde{x}, \tilde{y}). \tag{A.9}$$

with the prior factorizing and simplifying as

$$\begin{aligned} p(\theta, m, s \,|\, \tilde{x}, \tilde{y}) &= \tilde{p}(\theta_h \,|\, m, s, \tilde{x}, \tilde{y})\, p(\theta_L \,|\, m, s)\, p(m)\, p(s) \\ &\propto p(\theta_L \,|\, m, s)\, \tilde{p}(\tilde{y} \,|\, \tilde{x}, \theta; f)\, p(\theta_h)\, p(m)\, p(s), \end{aligned} \tag{A.10}$$

all of which we can compute analytically. With this prior, we can now derive a variational objective and perform approximate inference.

We begin by defining a variational distribution,

$$q(\theta, m, s, \tilde{x}, \tilde{y}) \doteq q(\theta_h)\, q(\theta_L \,|\, m, s)\, q(m)\, q(s)\, q(\tilde{x}, \tilde{y}),$$

and frame the inference problem of finding the posterior $p(\theta, m, s, \tilde{x}, \tilde{y} \,|\, x_\mathcal{D}, y_\mathcal{D})$ as a problem of optimization,

$$\min_{q_{\Theta, M, S, \tilde{X}, \tilde{Y}} \in \mathcal{Q}} D_{\text{KL}}(q_{\Theta, M, S, \tilde{X}, \tilde{Y}} \,\|\, p_{\Theta, M, S, \tilde{X}, \tilde{Y} \,|\, X_\mathcal{D}, Y_\mathcal{D}}),$$

where $\mathcal{Q}$ is a variational family. If the posterior $p_{\Theta, M, S, \tilde{X}, \tilde{Y} | X_\mathcal{D}, Y_\mathcal{D}}$ is in the variational family $\mathcal{Q}$, then the solution to the variational minimization problem is equal to the exact posterior. Modifying the inference problem by defining $q(\tilde{x}, \tilde{y}) \doteq p(\tilde{x}, \tilde{y}) = p(\tilde{y} \,|\, \tilde{x})p(\tilde{x})$, which further constrains the variational family, the optimization problem simplifies to

$$\min_{q_{\Theta, M, S} \in \mathcal{Q}} \mathbb{E}_{p_{\tilde{X}, \tilde{Y}}} \left[ D_{\text{KL}}(q_{\Theta, M, S} \,\|\, p_{\Theta, M, S \,|\, \tilde{X}, \tilde{Y}, X_\mathcal{D}, Y_\mathcal{D}}) \right],$$

which can equivalently be expressed as maximizing the variational objective

$$\bar{\mathcal{F}}(q_\Theta, q_M, q_S) \doteq \mathbb{E}_{q_{\Theta,M,S}}[\log p(y_\mathcal{D} \,|\, x_\mathcal{D}, \Theta; f)] - \mathbb{E}_{p_{\tilde{X},\tilde{Y}}}[D_{\mathrm{KL}}(q_{\Theta,M,S} \,\|\, p_{\Theta,M,S \,|\, \tilde{X},\tilde{Y}})].$$

To obtain a tractable expression of the regularization term, we first note that we can write

$$\mathbb{E}_{p_{\tilde{X},\tilde{Y}}}[D_{\mathrm{KL}}(q_{\Theta,M,S} \,\|\, p_{\Theta,M,S \,|\, \tilde{X},\tilde{Y}})]] = \mathbb{E}_{p_{\tilde{X},\tilde{Y}}}\Big[\mathbb{E}_{q_\Theta q_M q_S}[\log q(\Theta)q(M)q(S)] - \mathbb{E}_{q_\Theta q_M q_S}[\log p(\Theta, M, S \,|\, \tilde{X}, \tilde{Y})]\Big], \tag{A.11}$$

where the first term is the negative entropy and the second term is the negative cross-entropy. Using the same insights as above, we can write

$$\mathbb{E}_{p_{\tilde{X},\tilde{Y}}}[\mathbb{E}_{q_\Theta q_M q_S}[\log p(\Theta, M, S \,|\, \tilde{X}, \tilde{Y})]] \propto \mathbb{E}_{p_{\tilde{X},\tilde{Y}}}\Big[\mathbb{E}_{q_{\Theta_h} q_M q_S}[\log \tilde{p}(\tilde{Y} \,|\, \tilde{X}, \Theta_h, M, S)] \\ + \mathbb{E}_{q_\Theta}[\log p(\Theta_h)\, p(\Theta_L \,|\, M, S)\, p(M)\, p(S)]\Big], \tag{A.12}$$

up to an additive constant independent of $\theta$, and use it to express the KL divergence in Equation (A.11) up to an additive constant independent of $\theta$ as

$$D_{\mathrm{KL}}(q_{\Theta,M,S} \,\|\, p_{\Theta,M,S \,|\, \tilde{X},\tilde{Y}}) \propto -\mathbb{E}_{q_M q_S}[\mathbb{E}_{q_\Theta}[\log \tilde{p}(\tilde{Y} \,|\, \tilde{X}, \Theta_h, M, S)] \\ + D_{\mathrm{KL}}(q_{\Theta_L \,|\, M,S} \,\|\, p_{\Theta_L \,|\, M,S}) + D_{\mathrm{KL}}(q_{\Theta_h} \,\|\, p_{\Theta_h}) \\ + D_{\mathrm{KL}}(q_M \,\|\, p_M) + D_{\mathrm{KL}}(q_S \,\|\, p_S).$$

Now, further specifying the variational family as

$$\begin{aligned} q(\theta_L \,|\, m, s) &= \mathcal{N}(\theta_L; m, sI) \\ q(\theta_h) &= \mathcal{N}(\theta_h; \mu_h, \Sigma_h) \\ q(m) &= \mathcal{N}(m; \mu_L, \Sigma_L) \\ q(s) &= \mathrm{Lognormal}(s; \Sigma_L, \sigma_s^2 I). \end{aligned} \tag{A.13}$$

with learnable variational parameters $\mu \doteq \{\mu_h, \mu_m\}$ and $\Sigma \doteq \{\Sigma_L, \Sigma_L\}$ and fixed parameters $\{\sigma_m^2, \sigma_s^2\}$, we get $D_{\mathrm{KL}}(q_{\Theta_L \,|\, M,S} \,\|\, p_{\Theta_L \,|\, M,S}) = 0$, and the KL divergence simplifies to

$$D_{\mathrm{KL}}(q_{\Theta,M,S} \,\|\, p_{\Theta,M,S \,|\, \tilde{X},\tilde{Y}}) \propto -\mathbb{E}_{q_{\Theta_h} q_M q_S}[\log \tilde{p}(\tilde{Y} \,|\, \tilde{X}, \Theta_h, M, S)] + D_{\mathrm{KL}}(q_{\Theta_h} \,\|\, p_{\Theta_h}) \\ + D_{\mathrm{KL}}(q_M \,\|\, p_M) + D_{\mathrm{KL}}(q_S \,\|\, p_S), \tag{A.14}$$

where each of the KL divergences can be computed analytically, and we can obtain an unbiased estimator of the negative log-likelihood using simple Monte Carlo estimation.

Since $\Theta_h$ and $q_M$ are both mean-field Gaussian distributions, we can equivalently express the full variational objective in a simplified form as

$$\mathcal{F}(\mu, \Sigma) \doteq \underbrace{\mathbb{E}_{q_\Theta q_M q_S}[\log p(y_\mathcal{D} \,|\, x_\mathcal{D}, \Theta; f)]}_{\text{Expected log-likelihood}} - \underbrace{D_{\mathrm{KL}}(q_\Phi \,\|\, p_\Phi)}_{\text{KL regularization}} + \underbrace{\mathbb{E}_{q_{\Theta_h} q_M q_S}[\mathbb{E}_{p_{\tilde{X},\tilde{Y}}}[\log \tilde{p}(\tilde{Y} \,|\, \tilde{X}, \Theta_h, M, S)]] - \tau_s \|\Sigma_L\|_2^2}_{\text{Uncertainty regularization}}, \tag{A.15}$$

where we defined $\Phi \doteq \{\Theta_h, M\}$. We can estimate the expectations in the objective using simple Monte Carlo estimation, and gradients can be estimated using reparameterization gradients as in Blundell et al. (2015).

Letting $p_{\tilde{Y}|\tilde{X}}(\tilde{y} \,|\, \tilde{x}) = \delta(\mathbf{0})$ to encourage high uncertainty in the predictions on the set of context points, where $\delta(\cdot)$ is the Dirac delta function, we obtain the simplified objective

$$\mathcal{F}(\mu, \Sigma) \doteq \underbrace{\mathbb{E}_{q_\Theta q_M q_S}[\log p(y_\mathcal{D} \,|\, x_\mathcal{D}, \Theta; f)]}_{\text{Expected log-likelihood}} - \underbrace{D_{\mathrm{KL}}(q_\Phi \,\|\, p_\Phi)}_{\text{KL regularization}} + \underbrace{\mathbb{E}_{q_{\Theta_h} q_M q_S}[\mathbb{E}_{p_{\tilde{X}}}[\log \tilde{p}(\mathbf{0} \,|\, \tilde{X}, \Theta_h, M, S)]] - \tau_s \|\Sigma_L\|_2^2}_{\text{Uncertainty regularization}} \tag{A.16}$$

## Appendix B. Experimental details

### B.1. Training details

For model training, we use the joint fusion protocol defined by Hayat et al. (2022) in which the network is trained end-to-end including the modality specific encoders $\Phi_{cxr}$ and $\Phi_{ehr}$ using the fully connected layer $g(\cdot)$ to obtain the multi-label probabilities $\hat{y}_{\text{fusion}}$. Table A1 shows the details of the dataset splits used as input for our network.

Table A1: Summary of dataset sizes for the unimodal dataset and the combined multimodal dataset. We note that the size of the multimodal dataset decreases when the two modalities are paired.

| Dataset | Training | Validation | Testing | Context |
|---|---|---|---|---|
| Clinical time series data | 124,671 | 8,813 | 20,747 | 124,671 |
| Chest X-rays | 42,628 | 4,802 | 11,914 | 42,628 |
| Multimodal | 7,756 | 877 | 2,161 | 7,756 |

We use the binary cross-entropy loss (Good, 1952), adapted to the multi-label classification task:

$$\log p(y|x,\theta;f) = -\sum_{i=1}^{n}(y_i \log(\hat{y}_i) + (1-y_i)(\log(1-\hat{y}_i))), \tag{B.17}$$

where $\hat{y}_i \doteq \text{sigmoid}(f(x_i;\theta))$. The overall variational objective in our method is given by an expected log-likelihood term, KL regularization, and uncertainty regularization. In the stochastic setting, as described in Figure 1, we combine the training and context datasets as the input for the computation of this loss.

### B.2. Hyperparameter tuning

Initially, we used the deterministic baseline model to randomly sample a learning rate between $10^{-5}$ and $10^{-3}$ and selected the model and learning rate that achieved the model checkpoint with the best AUROC on the respective validation set. The best learning rate obtained was $2 \times 10^{-4}$, validated over 10 random seeds of training the deterministic model.

For the stochastic model, we performed a standard grid-based search to obtain the best hyper parameters for the regularization function. Table A2 shows the value ranges for each hyperparameter of our grid, which consists of 324 different model combinations. We note that this procedure requires more resources due to the higher number of hyperparameters as compared to the deterministic model. In addition, stochastic models also have more learnable parameters. In our case two times as many parameters, since the model has mean and variance parameters, and the regularization term requires performing a forward pass on the number of context points sampled from the context distribution (which we choose to be fewer points than are contained in each minibatch). In total, as is the case with every mean-field variational distribution, we have more learnable parameters than in a deterministic neural network and require more forward passes for every gradient step.

Table A2: Hyperparameter grid search values for the stochastic model

| Hyperparameter | Values | Best |
|---|---|---|
| prior variance | [1, 0.1, 0.01] | **0.1** |
| prior likelihood scale | [1, 0.1, 10] | **1** |
| prior likelihood f-scale | [0, 1, 10] | **10** |
| prior likelihood covariance scale | [0.1, 0.01, 0.001, 0.0001] | **0.1** |
| prior likelihood covariance diagonal | [1, 5, 0.5] | **5** |

### B.3. Model selection

We trained our stochastic model for 400 epochs. Since we have four metrics of interest (i.e., AUROC, AUPRC, Selective AUROC and Selective AUPRC), we computed the hypervolume using the volume formula of a 4-dimensional sphere as the main aggregated metric to select the best model checkpoint during training.

$$\text{hypervolume} = \frac{\pi^2 \text{R}^4}{2} \tag{B.18}$$

where R is the Euclidian magnitude of a 4-dimensional vector. The hypervolume approach ensures that we do not overfit to a single metric in the process of finding the best model.

### B.4. Technical implementation

Our data loading and pre-processing pipeline was implemented using PyTorch (Paszke et al., 2019) following the same structure of the code used by Hayat et al. (2022). However, we refactored the original unimodal and multimodal models, training, and evaluation loops using JAX (Bradbury et al., 2018). This framework simplifies the implementation of Bayesian neural networks and stochastic training, which are the basis of the uncertainty quantification methods used in this work. In addition, we obtained a significant reduction in total training time for the unimodal and multimodal models using JAX, compared to PyTorch.

We note that due to specific caching procedures of the JAX framework, we had to standardize each $x_{\text{ehr}}$ instance into 300 time steps for the LSTM encoder to avoid out-of-memory issues. The JAX framework requires that an LSTM encoder defines a static length of the sequences it is going to process, and then it caches this model in order to increase the training speed. This means that if different sequence lengths are used, then JAX would cache an instance of the LSTM encoder for each specific length to be used during each training cycle. The problem arises when dealing with a dataset that contains sequences of dynamic lengths that present high variance, i.e many different sequence lengths for every datapoint in the dataset, just as is the case with MIMIC-IV (Johnson et al., 2021). In comparison, PyTorch does not use this approach and is able to process sequences of dynamic length with one single instance of the LSTM encoder, however this is done at the cost of training speed when you compare both frameworks.

All of the experiments were executed using NVIDIA A100 and V100 80Gb Tensor Core GPUs.

## Appendix C. Additional Experimental Results

In this section, we provide additional results on the test set. Table A3 presents the results of the stochastic model for different context batch size values.

Table A3: Performance results on the test set for the stochastic model for varying values for context batch size.

| Context batch size | AUROC | AUPRC | Selective AUROC | Selective AUPRC |
|---|---|---|---|---|
| 16 | 0.732 (0.725, 0.739) | 0.511 (0.502, 0.525) | 0.740 (0.728, 0.753) | 0.447 (0.432, 0.469) |
| 32 | 0.733 (0.725, 0.739) | 0.510 (0.500, 0.524) | 0.743 (0.733, 0.756) | 0.448 (0.435, 0.466) |
| 64 | **0.735** (0.728, 0.742) | **0.514** (0.504, 0.528) | **0.748** (0.738, 0.760) | **0.452** (0.441, 0.472) |
| 128 | 0.733 (0.726, 0.739) | 0.512 (0.502, 0.525) | 0.728 (0.718, 0.739) | 0.401 (0.391, 0.418) |

Table A4, Table A5 and Table A6 present the extended results of our experiments for each label for the deterministic baseline, the Bayesian model with standard prior and the Bayesian model with M2D2 prior, respectively.

Table A4: Performance results across the different labels on the test set for the deterministic baseline (Hayat et al., 2022).

| Label | Prevalence | AUROC | AUPRC | Selective AUROC | Selective AUPRC |
|---|---|---|---|---|---|
| 1 Acute and unspecified renal failure | 0.321 | 0.753 (0.732, 0.774) | 0.574 (0.537, 0.613) | 0.775 (0.750, 0.802) | 0.655 (0.600, 0.717) |
| 2 Acute cerebrovascular disease | 0.078 | 0.854 (0.819, 0.888) | 0.427 (0.353, 0.512) | 0.595 (0.526, 0.687) | 0.072 (0.066, 0.098) |
| 3 Acute myocardial infarction | 0.093 | 0.692 (0.654, 0.727) | 0.189 (0.153, 0.234) | 0.634 (0.585, 0.675) | 0.089 (0.080, 0.107) |
| 4 Cardiac dysrhythmias | 0.379 | 0.690 (0.668, 0.713) | 0.563 (0.530, 0.599) | 0.678 (0.649, 0.705) | 0.621 (0.573, 0.682) |
| 5 Chronic kidney disease | 0.240 | 0.755 (0.732, 0.779) | 0.495 (0.454, 0.541) | 0.855 (0.830, 0.878) | 0.600 (0.529, 0.689) |
| 6 Chronic obstructive pulmonary disease | 0.148 | 0.735 (0.706, 0.765) | 0.338 (0.293, 0.390) | 0.814 (0.776, 0.850) | 0.431 (0.334, 0.536) |
| 7 Complications of surgical/medical care | 0.226 | 0.677 (0.650, 0.704) | 0.382 (0.339, 0.428) | 0.583 (0.552, 0.621) | 0.212 (0.202, 0.243) |
| 8 Conduction disorders | 0.115 | 0.787 (0.750, 0.822) | 0.575 (0.515, 0.636) | 0.849 (0.814, 0.882) | 0.746 (0.691, 0.800) |
| 9 Congestive heart failure; nonhypertensive | 0.295 | 0.772 (0.750, 0.794) | 0.593 (0.554, 0.632) | 0.829 (0.808, 0.853) | 0.705 (0.652, 0.758) |
| 10 Coronary atherosclerosis and related | 0.337 | 0.764 (0.744, 0.784) | 0.624 (0.583, 0.663) | 0.842 (0.814, 0.866) | 0.700 (0.640, 0.762) |
| 11 Diabetes mellitus with complications | 0.120 | 0.848 (0.823, 0.872) | 0.485 (0.417, 0.552) | 0.757 (0.704, 0.831) | 0.250 (0.205, 0.305) |
| 12 Diabetes mellitus without complication | 0.211 | 0.710 (0.684, 0.737) | 0.359 (0.324, 0.402) | 0.680 (0.645, 0.731) | 0.251 (0.219, 0.306) |
| 13 Disorders of lipid metabolism | 0.406 | 0.694 (0.671, 0.715) | 0.593 (0.556, 0.630) | 0.749 (0.721, 0.776) | 0.671 (0.617, 0.720) |
| 14 Essential hypertension | 0.433 | 0.653 (0.630, 0.676) | 0.561 (0.525, 0.595) | 0.624 (0.598, 0.653) | 0.599 (0.549, 0.653) |
| 15 Fluid and electrolyte disorders | 0.454 | 0.711 (0.689, 0.731) | 0.681 (0.649, 0.713) | 0.688 (0.664, 0.713) | 0.779 (0.746, 0.811) |
| 16 Gastrointestinal hemorrhage | 0.071 | 0.629 (0.583, 0.677) | 0.135 (0.100, 0.183) | 0.590 (0.542, 0.646) | 0.078 (0.066, 0.091) |
| 17 Hypertension with complications | 0.222 | 0.746 (0.720, 0.768) | 0.452 (0.407, 0.499) | 0.842 (0.810, 0.868) | 0.549 (0.451, 0.644) |
| 18 Other liver diseases | 0.169 | 0.684 (0.654, 0.715) | 0.336 (0.291, 0.385) | 0.701 (0.642, 0.754) | 0.375 (0.286, 0.476) |
| 19 Other lower respiratory disease | 0.126 | 0.615 (0.580, 0.651) | 0.209 (0.172, 0.256) | 0.577 (0.539, 0.612) | 0.118 (0.109, 0.141) |
| 20 Other upper respiratory disease | 0.054 | 0.638 (0.585, 0.686) | 0.092 (0.068, 0.127) | 0.551 (0.489, 0.640) | 0.055 (0.049, 0.074) |
| 21 Pleurisy; pneumothorax; pulmonary | 0.095 | 0.665 (0.629, 0.698) | 0.182 (0.146, 0.230) | 0.607 (0.558, 0.658) | 0.088 (0.084, 0.111) |
| 22 Pneumonia | 0.185 | 0.733 (0.707, 0.758) | 0.373 (0.327, 0.427) | 0.739 (0.698, 0.787) | 0.316 (0.261, 0.407) |
| 23 Respiratory failure; insufficiency; | 0.282 | 0.786 (0.766, 0.807) | 0.603 (0.566, 0.642) | 0.841 (0.817, 0.864) | 0.719 (0.671, 0.769) |
| 24 Septicemia (except in labor) | 0.227 | 0.755 (0.731, 0.778) | 0.504 (0.460, 0.550) | 0.841 (0.809, 0.869) | 0.626 (0.558, 0.696) |
| 25 Shock | 0.174 | 0.816 (0.791, 0.840) | 0.554 (0.507, 0.604) | 0.867 (0.821, 0.903) | 0.663 (0.580, 0.728) |

Table A5: Performance results across the different labels on the test set for the Bayesian (standard prior) model.

| Label | Prevalence | AUROC | AUPRC | Selective AUROC | Selective AUPRC |
|---|---|---|---|---|---|
| 1 Acute and unspecified renal failure | 0.321 | 0.747 (0.726, 0.768) | 0.573 (0.537, 0.611) | 0.817 (0.790, 0.845) | 0.672 (0.618, 0.732) |
| 2 Acute cerebrovascular disease | 0.078 | 0.861 (0.828, 0.893) | 0.418 (0.350, 0.498) | 0.590 (0.520, 0.678) | 0.074 (0.067, 0.101) |
| 3 Acute myocardial infarction | 0.093 | 0.705 (0.670, 0.741) | 0.208 (0.166, 0.266) | 0.694 (0.656, 0.729) | 0.087 (0.078, 0.104) |
| 4 Cardiac dysrhythmias | 0.379 | 0.698 (0.676, 0.721) | 0.577 (0.544, 0.615) | 0.744 (0.707, 0.776) | 0.626 (0.570, 0.688) |
| 5 Chronic kidney disease | 0.240 | 0.739 (0.716, 0.762) | 0.476 (0.432, 0.520) | 0.808 (0.772, 0.841) | 0.574 (0.499, 0.648) |
| 6 Chronic obstructive pulmonary disease | 0.148 | 0.729 (0.700, 0.761) | 0.338 (0.290, 0.393) | 0.802 (0.762, 0.842) | 0.428 (0.334, 0.535) |
| 7 Complications of surgical/medical care | 0.226 | 0.657 (0.627, 0.685) | 0.383 (0.338, 0.429) | 0.672 (0.627, 0.716) | 0.394 (0.321, 0.473) |
| 8 Conduction disorders | 0.115 | 0.817 (0.781, 0.848) | 0.596 (0.537, 0.657) | 0.871 (0.832, 0.906) | 0.734 (0.650, 0.810) |
| 9 Congestive heart failure; nonhypertensive | 0.295 | 0.783 (0.762, 0.805) | 0.618 (0.579, 0.656) | 0.867 (0.840, 0.892) | 0.739 (0.688, 0.787) |
| 10 Coronary atherosclerosis and related | 0.337 | 0.766 (0.745, 0.787) | 0.645 (0.605, 0.680) | 0.836 (0.802, 0.862) | 0.735 (0.686, 0.781) |
| 11 Diabetes mellitus with complications | 0.120 | 0.811 (0.783, 0.837) | 0.429 (0.368, 0.495) | 0.844 (0.804, 0.874) | 0.359 (0.240, 0.444) |
| 12 Diabetes mellitus without complication | 0.211 | 0.691 (0.664, 0.717) | 0.354 (0.319, 0.399) | 0.642 (0.607, 0.682) | 0.222 (0.201, 0.265) |
| 13 Disorders of lipid metabolism | 0.406 | 0.685 (0.662, 0.708) | 0.579 (0.545, 0.614) | 0.759 (0.730, 0.785) | 0.617 (0.564, 0.674) |
| 14 Essential hypertension | 0.433 | 0.660 (0.637, 0.681) | 0.576 (0.543, 0.611) | 0.716 (0.684, 0.746) | 0.632 (0.582, 0.680) |
| 15 Fluid and electrolyte disorders | 0.454 | 0.714 (0.693, 0.734) | 0.668 (0.635, 0.700) | 0.708 (0.685, 0.736) | 0.756 (0.718, 0.796) |
| 16 Gastrointestinal hemorrhage | 0.071 | 0.638 (0.593, 0.682) | 0.131 (0.097, 0.177) | 0.606 (0.561, 0.676) | 0.071 (0.066, 0.092) |
| 17 Hypertension with complications | 0.222 | 0.733 (0.710, 0.757) | 0.431 (0.388, 0.479) | 0.779 (0.735, 0.818) | 0.486 (0.397, 0.579) |
| 18 Other liver diseases | 0.169 | 0.696 (0.664, 0.727) | 0.354 (0.305, 0.403) | 0.731 (0.680, 0.780) | 0.429 (0.345, 0.526) |
| 19 Other lower respiratory disease | 0.126 | 0.604 (0.567, 0.642) | 0.181 (0.153, 0.219) | 0.582 (0.547, 0.617) | 0.123 (0.113, 0.145) |
| 20 Other upper respiratory disease | 0.054 | 0.685 (0.632, 0.739) | 0.165 (0.109, 0.236) | 0.618 (0.552, 0.685) | 0.101 (0.052, 0.205) |
| 21 Pleurisy; pneumothorax; pulmonary | 0.095 | 0.666 (0.629, 0.701) | 0.166 (0.135, 0.208) | 0.593 (0.543, 0.657) | 0.090 (0.082, 0.120) |
| 22 Pneumonia | 0.185 | 0.758 (0.732, 0.781) | 0.400 (0.355, 0.455) | 0.783 (0.743, 0.829) | 0.341 (0.270, 0.436) |
| 23 Respiratory failure; insufficiency; | 0.282 | 0.824 (0.804, 0.843) | 0.631 (0.592, 0.675) | 0.890 (0.868, 0.909) | 0.703 (0.637, 0.771) |
| 24 Septicemia (except in labor) | 0.227 | 0.783 (0.761, 0.805) | 0.522 (0.476, 0.572) | 0.866 (0.834, 0.892) | 0.629 (0.530, 0.703) |
| 25 Shock | 0.174 | 0.826 (0.804, 0.847) | 0.552 (0.502, 0.606) | 0.888 (0.858, 0.918) | 0.582 (0.507, 0.658) |

Table A6: Performance results across the different labels on the test set for the Bayesian (M2D2 prior) model.

| Label | Prevalence | AUROC | AUPRC | Selective AUROC | Selective AUPRC |
|---|---|---|---|---|---|
| 1 Acute and unspecified renal failure | 0.321 | 0.756 (0.735, 0.779) | 0.587 (0.551, 0.627) | 0.830 (0.800, 0.857) | 0.680 (0.617, 0.740) |
| 2 Acute cerebrovascular disease | 0.078 | 0.870 (0.840, 0.901) | 0.459 (0.385, 0.545) | 0.664 (0.598, 0.745) | 0.088 (0.076, 0.115) |
| 3 Acute myocardial infarction | 0.093 | 0.716 (0.681, 0.754) | 0.220 (0.174, 0.277) | 0.656 (0.611, 0.698) | 0.083 (0.077, 0.104) |
| 4 Cardiac dysrhythmias | 0.379 | 0.687 (0.663, 0.710) | 0.570 (0.536, 0.606) | 0.737 (0.706, 0.768) | 0.654 (0.603, 0.709) |
| 5 Chronic kidney disease | 0.240 | 0.767 (0.744, 0.789) | 0.507 (0.464, 0.553) | 0.853 (0.825, 0.879) | 0.612 (0.534, 0.690) |
| 6 Chronic obstructive pulmonary disease | 0.148 | 0.727 (0.700, 0.757) | 0.326 (0.280, 0.377) | 0.773 (0.725, 0.814) | 0.413 (0.293, 0.503) |
| 7 Complications of surgical/medical care | 0.226 | 0.659 (0.631, 0.686) | 0.396 (0.351, 0.444) | 0.666 (0.620, 0.708) | 0.375 (0.310, 0.462) |
| 8 Conduction disorders | 0.115 | 0.798 (0.763, 0.832) | 0.593 (0.532, 0.656) | 0.850 (0.814, 0.889) | 0.754 (0.689, 0.812) |
| 9 Congestive heart failure; nonhypertensive | 0.295 | 0.788 (0.768, 0.808) | 0.600 (0.562, 0.637) | 0.874 (0.853, 0.892) | 0.722 (0.663, 0.773) |
| 10 Coronary atherosclerosis and related | 0.337 | 0.767 (0.746, 0.787) | 0.626 (0.588, 0.665) | 0.846 (0.820, 0.869) | 0.716 (0.654, 0.769) |
| 11 Diabetes mellitus with complications | 0.120 | 0.842 (0.817, 0.866) | 0.461 (0.398, 0.528) | 0.821 (0.765, 0.906) | 0.344 (0.256, 0.472) |
| 12 Diabetes mellitus without complication | 0.211 | 0.716 (0.690, 0.742) | 0.386 (0.345, 0.432) | 0.673 (0.634, 0.727) | 0.256 (0.223, 0.313) |
| 13 Disorders of lipid metabolism | 0.406 | 0.698 (0.675, 0.720) | 0.591 (0.558, 0.628) | 0.773 (0.746, 0.798) | 0.624 (0.572, 0.685) |
| 14 Essential hypertension | 0.433 | 0.669 (0.646, 0.694) | 0.588 (0.555, 0.621) | 0.709 (0.677, 0.740) | 0.643 (0.588, 0.695) |
| 15 Fluid and electrolyte disorders | 0.454 | 0.717 (0.695, 0.737) | 0.679 (0.647, 0.709) | 0.739 (0.712, 0.767) | 0.771 (0.732, 0.811) |
| 16 Gastrointestinal hemorrhage | 0.071 | 0.667 (0.627, 0.708) | 0.134 (0.104, 0.185) | 0.596 (0.537, 0.720) | 0.068 (0.064, 0.089) |
| 17 Hypertension with complications | 0.222 | 0.760 (0.736, 0.781) | 0.475 (0.430, 0.522) | 0.843 (0.810, 0.872) | 0.574 (0.483, 0.666) |
| 18 Other liver diseases | 0.169 | 0.723 (0.693, 0.750) | 0.378 (0.330, 0.428) | 0.673 (0.614, 0.743) | 0.313 (0.228, 0.423) |
| 19 Other lower respiratory disease | 0.126 | 0.594 (0.560, 0.629) | 0.181 (0.154, 0.222) | 0.568 (0.533, 0.612) | 0.120 (0.114, 0.144) |
| 20 Other upper respiratory disease | 0.054 | 0.670 (0.614, 0.722) | 0.133 (0.094, 0.193) | 0.540 (0.483, 0.603) | 0.052 (0.048, 0.071) |
| 21 Pleurisy; pneumothorax; pulmonary | 0.095 | 0.692 (0.658, 0.726) | 0.167 (0.139, 0.207) | 0.630 (0.576, 0.707) | 0.086 (0.081, 0.109) |
| 22 Pneumonia | 0.185 | 0.756 (0.732, 0.781) | 0.411 (0.362, 0.468) | 0.808 (0.769, 0.851) | 0.395 (0.328, 0.499) |
| 23 Respiratory failure; insufficiency; | 0.282 | 0.811 (0.790, 0.831) | 0.633 (0.594, 0.673) | 0.876 (0.848, 0.899) | 0.733 (0.675, 0.791) |
| 24 Septicemia (except in labor) | 0.227 | 0.774 (0.751, 0.797) | 0.513 (0.467, 0.559) | 0.818 (0.772, 0.865) | 0.582 (0.504, 0.664) |
| 25 Shock | 0.174 | 0.809 (0.786, 0.833) | 0.536 (0.486, 0.586) | 0.876 (0.838, 0.902) | 0.637 (0.546, 0.702) |