CALIBRATING MULTIMODAL LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Multimodal machine learning has achieved remarkable progress in a wide range of scenarios. However, the reliability of multimodal learning remains largely unexplored. In this paper, through extensive empirical studies, we identify current methods suffer from unreliable predictive confidence that tends to rely on partial modalities when estimating confidence. Specifically, we find that the confidence estimated by current models could even increase when some modalities are corrupted. To address the issue, we introduce an intuitive principle for multimodal classification, i.e., the confidence should not increase when one modality is removed. Accordingly, we propose a novel regularization technique, i.e., Calibrating Multimodal Learning (CML) regularization, to calibrate the predictive confidence of previous methods. This technique could be flexibly equipped by existing models and improve the performance in terms of confidence calibration, classification accuracy, and model robustness.

1 INTRODUCTION

Multiple modalities data widely exist in real-world applications such as medical analysis (Perrin et al., 2009), social media (Wang et al., 2019), and autonomous driving (Khodayari et al., 2010). To fully explore the potential value of each modality, multimodal learning emerges as a promising way to train a machine learning (ML) model by integrating all available multimodal cues for further data analysis tasks (e.g., instance classification). Numerous approaches have been proposed to build multimodal classification paradigms for various tasks (Wang et al., 2019; Antol et al., 2015; Bagher Zadeh et al., 2018; Kishi et al., 2019). Despite above progresses, the reliability of current multimodal classification methods remains largely unexplored. One key aspect of the reliability is to build a high-quality uncertainty estimator (Neal, 2012; MacKay, 1992), which can quantitatively characterize the probability that predictions will be wrong. With such an estimator, further processing can be taken to improve the performance of the system (e.g., human assistance) when the predictive uncertainty is high. This is especially useful in high-stake scenarios (Hafner et al., 2019; Qaddoum & Hines, 2012).

In the setting of multimodal classification, in addition to exact overall prediction confidence, the relationship between the modalities should also be taken into concerns. Intuitively, the confidence of an ideal multimodal classifier should not increase when one modality is removed. An illustrative example of an ideal confidence estimator is shown in Fig. 1, where the confidence gradually decreases when the observed information becomes less comprehensive. However, in practice, the confidence estimation obtained by ordinary training methods tends to be overconfident to partial modalities, which violates the principle and leads to some counter-intuitive phenomenons. We perform extensive empirical studies and observe that when a modality is removed, the overall confidence can even increase. This observation contradicts the usual assumption of multimodal classification since modalities are assumed to be predictive of the target for most multimodal classification tasks (Wu et al., 2022). Intuitively, this implies that the models are more inclined to believe in a unique modality and is prone to be affected by this modality, which has also been shown in prior works (Wu et al., 2022; Wang et al., 2020). This further impairs the robustness of the learned models, i.e., the models are easy to be influenced when some modalities are corrupted, since the models can not make decisions taking all modalities into account fairly.

A natural idea to address the above issue is to employ recent uncertainty calibration methods such as temperature scaling (Guo et al., 2017) or Bayesian learning (Cobb & Jalaian, 2021; Karaletsos & Bui, 2020; Foong et al., 2020), which can build more accurate uncertainty estimation than the



Figure 1: Motivation of calibrating multimodal learning. The confidence of an ideal multimodal classifier should decrease (at least not increase) when one modality is removed.

traditional training/inference manner. However, these approaches do not explicitly consider the relationship between different modalities (i.e., they can only calibrate the fused confidence but cannot adjust the relationship between different modalities during training) and thus still fail to achieve satisfactory performance in the multimodal classification setting. Therefore, we claim that high-quality uncertainty estimation in various multimodal classification tasks needs to explicitly treat all modalities fairly. To this end, we propose a novel regularization technique called Calibrating Multimodal Learning (CML) which enforces the consistency between prediction confidence and the number of modalities. The motivation of CML is based on a natural intuition, i.e., the prediction confidence should decrease (at least not increase) when one modality is removed, which could intrinsically improve the confidence calibration. Specifically, we propose a simple regularization term that enforces a model to learn an intuitive ranking relationship by adding a penalty for the samples whose predictive confidence will increase when one modality is removed. The main contributions of this paper are summarized as follows:

- We perform extensive empirical studies to show that most existing multimodal classification paradigms tend to over-rely on partial modalities (different samples over-rely on different modalities rather than all samples over-rely on the same modalities), which implies they fail to achieve trustworthy uncertainty estimation.
- We introduce a measure to evaluate the reliability of the confidence estimation from the confidence ranking perspective, which can characterize whether a multimodal learning method can treat all modalities fairly.
- We propose a regularization strategy to calibrate the confidence of various multimodal classification methods. We then conduct extensive experiments to show the superiority of our method in terms of the confidence calibration (Table 1), classification accuracy (Table 2) and model robustness (Table 3).

2 Related work

Uncertainty estimation provides a way for trustworthy prediction (Abdar et al., 2021). Uncertainty can be used as an indicator of whether the predictions given by models are prone to be wrong. Many uncertainty-based models have been proposed in the past decades, such as Bayesian



Figure 2: Current methods (MMTM (Wu et al., 2022), CPM-Nets (Zhang et al., 2019), and MI-WAE (Mattei & Frellsen, 2019)) violate the Proposition 1 (red color indicates the proportion of test samples whose predictive confidence given by the model decreases while providing more modalities, "CI" is defined in Eq. 1). We estimate the performance on two-modality datasets, and the pie charts show that different samples over-rely on different modalities rather than all samples over-rely on the same modality (e.g., "53% Mod1" indicates "among the samples who violate Proposition 1, there is 53 percent of samples whose confidence will increase when Mod2 is removed and the other samples will increase confidence when Mod1 is removed").

neural networks (Neal, 2012; MacKay, 1992; Denker & LeCun, 1990; Kendall & Gal, 2017), Dropout (Molchanov et al., 2017), Deep ensembles (Lakshminarayanan et al., 2017; Havasi et al., 2020), and DUQ (van Amersfoort et al., 2020) built upon RBF networks. **Prediction confidence** is always referred to in classification models, which expects the predicted class probability to be consistent with the empirical accuracy. Many methods focus on smoothing the prediction probabilities distribution, such as Label smoothing (Müller et al., 2019), focal loss (Mukhoti et al., 2020), TCP (Corbière et al., 2019)and Temperature scaling (TS) (Guo et al., 2017). More related researches please refer to Appendix G.

Multimodal learning emerges as a promising way to train a machine learning model. Recently, there have been a wide range of research interests in handling missing modalities for multimodal learning, including imputation-independent (Type I) methods (Zhang et al., 2019) and imputation-dependent (Type II) methods (Mattei & Frellsen, 2019; Wu & Goodman, 2018). Imputation-independent methods have no need to reconstruct the missing modalities and make classification via an uniform representation. For imputation-dependent methods (based on reconstruction), the strategy model can be split into two stages, reconstructing the missing modalities and making classification according to the reconstructed modalities. Besides the methods for incomplete multimodal learning, recent multimodal methods (Type III) Joze et al. (2020); Wu et al. (2022) achieve SOTA performance in multimodal video classification. We evaluate the performance of CPM-Nets (Zhang et al., 2019), MIWAE (Mattei & Frellsen, 2019), and MMTM (Joze et al., 2020; Wu et al., 2022) in experiments due to their representativeness in different types of multimodal learning.

3 Method

In this section, we first introduce some basic notations in Section 3.1. Then, we show the basic assumption of our method and its empirical motivations in Section 3.2 and evaluate the confidence estimation performance of current multimodal methods in Section 3.3. At the end, we propose a simple yet effective regularization technique and elaborate the technical details in Section 3.4.

3.1 NOTATION

We define the training data as $\mathcal{D} = \left\{ \{x_i^m\}_{m=1}^M, y_i\}_{i=1}^N$, where x_i^m is the input feature of the *m*-th modality of the *i*-th sample, and $y_i \in \{1, \dots, K\}$ is the corresponding class label. To distinguish whether the input is a unique modality or a set of modalities, we use x^m to represent the *m*-th modality, and use the $\mathbf{x}^{(S)}$ to represent multimodal input set, where S is a set of modalities' indexes (e.g., if we have $S = \{1, 2\}$, then $\mathbf{x}^{(S)}$ indicates a feature set consisting of x^1 and x^2 , and $\mathbf{x}^{(M)} = \{\mathbf{x}^1, \dots, \mathbf{x}^M\}$ indicates the complete M modalities). The goal is to learn a function: $f(\mathbf{x}^{(M)}) \to z$, where the output z of the network is a vector of K values called logits. Then the logits vector is

transformed by a softmax layer: $\hat{p}_k = e^{z_k} / \sum_k e^{z_k}$, where the probability distribution of a sample x is defined as $P(y \mid \mathbf{w}, \mathbf{x}^{(\mathbb{M})}) = \{\hat{p}_k\}_1^K$. The predicted class label is: $\hat{y} = \arg \max_y P(y \mid \mathbf{w}, \mathbf{x}^{(\mathbb{M})})$.

3.2 BASIC ASSUMPTION

In the real-world applications, the quality of test samples is unstable (e.g., some modalities may be corrupted), so the quality (i.e., available modalities) of the multimodal input should be reflected in some quantitative manner (i.e., predictive confidence) when multimodal methods are deployed in the real-world applications. However, it is difficult to exactly define the "quality" of each sample, and we cannot define the exact functional relationship between the quality and confidence since the confidence estimation is unique for different models even for a same sample. To address this issue, we approximate this relationship with a ranking-based form as follow:

Proposition 1. Given two versions of a unique sample $\mathbf{x}^{(\mathbb{M})}$, i.e., $\mathbf{x}^{(\mathbb{T})}$ and $\mathbf{x}^{(\mathbb{S})}$, if we can assure $\mathbb{T} \subset \mathbb{S}$, then, for a trustworthy multimodal classifier $f(\cdot)$, it should hold $\operatorname{Conf}(f(\mathbf{x}^{(\mathbb{T})})) \leq \operatorname{Conf}(f(\mathbf{x}^{(\mathbb{S})})$.

Proposition 1 indicates that the predictive confidence shouldn't increase when one modality is removed. We further define the prediction Confidence Increment (CI) with informativeness increment for a unique sample as:

$$\operatorname{CI}(\mathbf{x}^{(\mathbb{S})}, \mathbf{x}^{(\mathbb{T})}) = \operatorname{Conf}(f(\mathbf{x}^{(\mathbb{S})})) - \operatorname{Conf}(f(\mathbf{x}^{(\mathbb{T})})) \quad \text{subject to: } \mathbb{T} \subset \mathbb{S},$$
(1)

where \mathbb{T} and \mathbb{S} are sets of modalities' indexes. Specially, a negative value indicates a poor confidence estimation performance that the predictive confidence increases when one modality is removed. To quantify the extent that a learned model violates Proposition 1, we introduce a novel measure: Violating Ranking Rate (VRR) as the proportion of test samples whose predictive confidence will increase when removing one modality:

$$\operatorname{VRR} = \mathbb{E}\left[\mathbb{1}\left(\operatorname{CI}(\mathbf{x}^{(\mathbb{S})}, \mathbf{x}^{(\mathbb{T})}) < 0\right)\right] \quad \text{subject to: } \mathbb{T} \subset \mathbb{S}.$$
(2)

We initialize S as the complete modalities, and obtain T by randomly removing a modality from S. Then T is regarded as S for another confidence ranking pair and we repeat this process until there is only one modality remained in T (please refer to Appendix A for detail). A natural question then arises: How about the confidence estimation performance of the current methods when one modality is removed?

3.3 CONFIDENCE ESTIMATION PERFORMANCE OF CURRENT MULTIMODAL METHODS

To evaluate the quality of confidence estimation of existing multimodal classifiers, we compute the VRR score of CPM-Nets (Zhang et al., 2019) and MIWAE (Mattei & Frellsen, 2019), which are two typical methods in handling incomplete multimodal data. In addition to incomplete multimodal, we also evaluate the MMTM (Wu et al., 2022), which is a SOTA multimodal classification method. As shown in Tab. 1, the VRR scores of previous methods are quite high which indicates the prediction confidence on many samples will violate Proposition 1. And the visualization is shown in Fig. 2, where the red color indicates the proportion of test samples whose predictive confidence given by the model decreases while providing more modalities.

A naive strategy is to re-balance the contribution of every modality (i.e., allocating a smaller weight to the modality that samples over-rely on during the fusion). As shown in Fig. 2, however, we find that different samples over-rely on different modalities rather than all samples over-rely on the same modality. This indicates that the problem can't be solved by re-weighting the overall contribution of different modalities since it will make the confidence estimation of some samples worse. Instead, our method characterizes the relationship between the modalities in term of simplewise, which inherently calibrates the contribution for all samples. Intuitively, it is risky for a model which usually increases the prediction confidence when one modality is removed, which implies that the confidence and its quality are not matched. For this issue, they cannot be deployed into risk-sensitive applications such as medical analysis. As a comparison, our method can significantly decrease VRR score (see more details in Tab. 1) for a more trustworthy confidence estimation.

3.4 CML REGULARIZATION

As shown in Section 3.3, the current multimodal methods usually increase the prediction confidence when one modality is removed, which hinders the model performance inherently. To address this issue, a naive strategy is to penalize the confidence difference between the $\mathbf{x}^{(\mathbb{T})}$ and $\mathbf{x}^{(\mathbb{S})}$:

$$\mathcal{L}^{(\mathbb{T},\mathbb{S})} = \operatorname{Conf}(\mathbf{x}^{(\mathbb{T})}) - \operatorname{Conf}(\mathbf{x}^{(\mathbb{S})}).$$
(3)

However, models sometimes can still make an accurate prediction confidently when one modality is removed in practice. Eq. 3 forces the models to predict small confidence strictly when one modality is removed, which pushes the model to estimate an extremely small confidence for each modality (the illustration please refer to Appendix C.6). For this issue, we relax this regularization by only penalizing the situation that the estimated confidence increases when one modality is removed. For any pair of multimodal inputs which satisfies that $\mathbb{T} \subset \mathbb{S}$, the regularization can be written as:

$$\mathcal{L}^{(\mathbb{T},\mathbb{S})} = \max\left(0, \operatorname{Conf}(\mathbf{x}^{(\mathbb{T})}) - \operatorname{Conf}(\mathbf{x}^{(\mathbb{S})})\right).$$
(4)

For each sample, the total regularization loss is integrated over all pairs of inputs with different numbers of modalities, which can be written as:

$$\mathcal{L}^{\text{CML}} = \sum_{\mathbb{T}, \,\mathbb{S}} \mathcal{L}^{(\mathbb{T}, \mathbb{S})}, \quad \{\forall (\mathbb{T}, \,\mathbb{S}) | \mathbb{T} \subset \mathbb{S}\}.$$
(5)

The exact computation of above loss needs to enumerate all modality set pairs (i.e., \mathbb{T} and \mathbb{S}), which is typically computational expensive sometimes. Therefore, we propose to approximate this loss by sampling modality set pairs and find this strategy works well in practice. Specifically, we conduct sampling as same as that in computing VRR (Eq. 2).

The proposed regularization is general and thus can be equipped by current multimodal classifiers to calibrate their confidence estimation as an additional loss item. We typically provide examples in utilizing the proposed technique in imputation-independent method(i.e., CPM-Nets (Zhang et al., 2019)), imputation-dependent method(i.e., MIWAE (Mattei & Frellsen, 2019)), and recent multimodal classification method (i.e., MMTM (Wu et al., 2022)). The proposed regularization can be deployed to current multimodal methods flexibly, the objective function is induced as:

$$\mathcal{L}_i = \mathcal{L}_i^{\rm cl} + \lambda \mathcal{L}_i^{\rm CML},\tag{6}$$

where \mathcal{L}_i^{cl} is the classification loss criterion (e.g., cross-entropy loss), and λ is the hyperparameter controlling the strength of CML regularization. The details are shown in Algorithm 1.

Algorithm 1: The training pseudocode for deploying CML regularization

1 Given dataset D = {{x_i^m}_{m=1}^M, y_i}^N_{i=1}, initialized classifier f, classification loss criterion L^{cl}, coefficient of CML regularization λ, and epochs for training the classifier train_epochs
 2 for e = 1,..., train_epochs do

 $\mathbb{S} \leftarrow \mathbb{M}; \ \mathcal{L}^{cl} \leftarrow \mathcal{L}^{cl}(\mathbf{x}^{(\mathbb{S})}); \ \mathcal{L}^{CML} \leftarrow 0$ 3 for m = 1, ..., M - 1 do 4 Randomly remove a modality of \mathbb{S} and set it as \mathbb{T} 5 Compute the classification loss: $\mathcal{L}^{cl} \leftarrow \mathcal{L}^{cl} + \mathcal{L}^{cl}(\mathbf{x}^{(\mathbb{T})})$ 6 Compute the regularization loss: 7 $\mathcal{L}^{\text{CML}} \leftarrow \mathcal{L}^{\text{CML}} + \max\left(0, \operatorname{Conf}(\mathbf{x}^{(\mathbb{T})}) - \operatorname{Conf}(\mathbf{x}^{(\mathbb{S})}) - \tau\right)$ $\mathbb{S} \leftarrow \mathbb{T}$ 8 end Total loss: $\mathcal{L} = \frac{1}{M} \mathcal{L}^{cl} + \lambda \mathcal{L}^{CML}$ 10 Update the parameters of the classifier f with \mathcal{L} 11 12 end 13 **return** the classifier f

3.5 DISCUSSION AND ANALYSES

• Why should a model meet the ranking relationship regardless of the label? For multimodal learning, all modalities are assumed to be predictive of the target (Wu et al., 2022), which can be expressed as $I(y, \mathbf{x}^m) \ge 0$, where $I(\cdot)$ denotes mutual information (Blum & Mitchell, 1998).

Lemma 3.1. Suppose we have two versions of a unique sample $\mathbf{x}^{(\mathbb{M})}$, i.e., $\mathbf{x}^{(\mathbb{T})}$ and $\mathbf{x}^{(\mathbb{S})}$, if we can assure $\mathbb{T} \subset \mathbb{S}$, then, for any class label y, we have $I(y, \mathbf{x}^{(\mathbb{T})}) \leq I(y, \mathbf{x}^{(\mathbb{S})})$.

In other words, $\mathbf{x}^{(S)}$ is more predictive of the target than $\mathbf{x}^{(T)}$ regardless of the label. For a trustworthy multimodal method, the confidence of $\mathbf{x}^{(T)}$ should not be larger than $\mathbf{x}^{(S)}$.

 \circ Why can CML regularization calibrate the model? CML regularization can guarantee a smaller confidence of $\mathbf{x}^{(\mathbb{T})}$ when the model makes a wrong prediction of $\mathbf{x}^{(\mathbb{S})}$, which means that CML can alleviate the over-confidence.

Lemma 3.2. Suppose the CML regularization can achieve a lower VRR, i.e., $\operatorname{VRR}_{CML} < \operatorname{VRR}_{Ori}$, then for the samples that meet $\mathbb{E}\left(\operatorname{Conf}_{CML}(\mathbf{x}^{(\mathbb{S})})\right) = \mathbb{E}\left(\operatorname{Conf}_{Ori}(\mathbf{x}^{(\mathbb{S})})\right)$, we have $\mathbb{E}\left(\operatorname{Conf}_{CML}(\mathbf{x}^{(\mathbb{T})})\right) \leq \mathbb{E}\left(\operatorname{Conf}_{Ori}(\mathbf{x}^{(\mathbb{T})})\right)$.

Although it is difficult to make $\operatorname{Conf}_{Ori}(\cdot)$ equal to $\operatorname{Conf}_{CML}(\cdot)$ strictly for all samples, as shown in Appendix C.5, we find $\operatorname{Conf}_{CML}(\mathbf{x}^{(\mathbb{S})})$ and $\operatorname{Conf}_{Ori}(\mathbf{x}^{(\mathbb{S})})$ are very similar for most samples, where $\operatorname{Conf}_{Ori}(\cdot)$ and $\operatorname{Conf}_{CML}(\cdot)$ indicate the confidence estimated by the original model and the model improved by CML regularization respectively.

 \circ Why not just penalize the difference in confidence: $\operatorname{Conf}(\mathbf{x}^{(\mathbb{T})}) - \operatorname{Conf}(\mathbf{x}^{(\mathbb{S})})$? Forcing the confidence for $\mathbf{x}^{(\mathbb{T})}$ to be smaller than the confidence for $\mathbf{x}^{(\mathbb{S})}$ strictly will lead to a very small confidence for $\mathbf{x}^{(\mathbb{T})}$ and will make the model estimate an extremely small confidence for each modality, and the experiments are shown in Appendix C.6. What's more, the model sometimes can still make correct predictions confidently when one modality is removed. A flexible ranking regularization makes it more suitable for real data.

4 EXPERIMENTS

4.1 Setup

We deploy the proposed regularization method to different types of multimodal classifiers including the imputation-independent method (Type I), the imputation-dependent method (Type II), and the recent SOTA method (Type III). CPM-Nets (Zhang et al., 2019) is a typical imputation-independent algorithm, which can adapt to arbitrary missing patterns without reconstructing the missing modalities. MIWAE (Mattei & Frellsen, 2019) is a typical imputation-dependent algorithm. The above two methods are typical in incomplete multimodal learning. In addition to incomplete multimodal learning methods, we also deploy the regularization to the current multimodal method (Wu et al., 2022), which is named MMTM. For MMTM, we approximate removing one modality by corrupting its features (e.g., adding strong noise) due to the model can't make a prediction when one modality is absolutely removed.

Datasets: We evaluate the proposed method on diverse datasets, including YaleB (Georghiades et al., 2002), Handwritten (Perkins & Theiler, 2003), CUB (Wah et al., 2011), Animal (Krizhevsky et al., 2012; Simonyan & Zisserman, 2014), TUANDROMD (Borah et al., 2020), NYUD2 (Qi et al., 2017), and SUNRGBD (Song et al., 2015).

Experiment setting: For a fair comparison, the only difference between whether the model is equipped with CML regularization or not is whether the coefficient λ is set to 0. Please refer to Appendix C.2 for more detailed settings.

4.2 QUESTIONS TO BE VERIFIED

We conduct diverse experiments to comprehensively investigate the underlying assumption and the proposed method, including:

• **Can CML regularization improve the confidence estimation of multimodal classifiers?** To validate whether the proposed method improves multimodal classifiers' confidence estimation, we evaluate the confidence estimation of current multimodal classifiers without and with CML regularization, respectively. We conduct experiments of each type of methods on five datasets and evaluate their trustworthiness in terms of VRR (defined in the Eq. 2).

• **Can CML regularization improve robustness?** CML regularization can improve multimodal classifiers' confidence estimation, so a natural question arises - does a better confidence estimation imply better robustness? To verify this, we evaluate the robustness on the complete multimodal data and noisy multimodal data (adding Gaussian noise to some modalities, i.e., zero mean with varying variance ϵ).

 \circ Is CML easy to be deployed and not sensitive to hyperparameters? In order to investigate the key factor that makes the improvement in the proposed method, we evaluate the performance in terms of classification accuracy under different strengths of CML regularization. We conduct experiments on both the original and noised data (i.e., adding noise to one of the modalities during the test). More details are shown in Appendix C.2.

Table 1: VRR (%) of test samples (a lower value indicates a better confidence estimation. Type III is shown in Appendix). " λ " indicates the model is not equipped with the proposed regularization ($\lambda = 0$). Performance on Type III please refer to Tab. 8 (Appendix F).

Method	CML	TUANDROMD	YaleB	Handwritten	CUB	Animal
Type I	x ✓ Improve	$\begin{array}{c} 23.38 \pm 1.39 \\ 12.58 \pm 2.84 \\ \bigtriangleup 10.80 \end{array}$	$\begin{array}{c} 39.15 \pm 4.97 \\ 15.05 \pm 1.12 \\ \bigtriangleup 24.10 \end{array}$	$\begin{array}{c} 17.64 \pm 2.31 \\ 3.18 \pm 0.80 \\ \bigtriangleup 14.46 \end{array}$	$2.83 \pm 1.55 \\ 2.17 \pm 1.13 \\ riangle 0.66$	$\begin{array}{c} 44.39 \pm 7.55 \\ 29.02 \pm 5.43 \\ \bigtriangleup 15.37 \end{array}$
Type II	× ✓ Improve	$\begin{array}{c} 39.17 \pm 2.32 \\ 8.38 \pm 1.31 \\ \bigtriangleup 30.79 \end{array}$	$\begin{array}{c} 20.54 \pm 4.26 \\ 14.46 \pm 2.17 \\ \bigtriangleup 6.08 \end{array}$	$\begin{array}{c} 33.82 \pm 5.16 \\ 29.99 \pm 2.30 \\ \bigtriangleup 3.83 \end{array}$	$\begin{array}{c} 23.17 \pm 4.87 \\ 20.17 \pm 3.05 \\ \bigtriangleup 3.00 \end{array}$	$\begin{array}{c} 12.51 \pm 1.50 \\ 8.64 \pm 0.32 \\ \bigtriangleup 3.87 \end{array}$

4.3 RESULTS

4.3.1 CONFIDENCE ESTIMATION

We evaluate the confidence estimation of current multimodal classification models from a ranking perspective and find that for a large number of samples the confidence will increase when one modality is removed, while the confidence estimation of classification models equipped with the proposed CML regularization is significantly improved. We intuitively demonstrate the confidence changing in Fig. 3, and the quantitative results are shown in Tab. 1. According to Fig. 3, we show the confidence estimation of CPM-Nets, where "Original" and "CML" indicate the model is without and with the proposed CML regularization respectively. Ac-



Figure 3: Confidence estimation when one modality is removed, where "CI" is defined in Eq. 1.

cording to Fig. 3, it is observed that the confidence without CML regularization may increase when one modality is removed, which indicates that the model doesn't take all modalities into account fairly when making predictions. This will lead to unpromising robustness and generalization, which clearly verifies the main assumption in Sec. 4.3.2.

We also report the quantitative results on five datasets. It is observed that the confidence estimation of each model is obviously improved with the proposed CML regularization.

Method	Dataset	CML	Accuracy (†)	$\begin{array}{c} \text{NLL} \\ (\downarrow) \end{array}$	AURC (↓)	E-AURC (↓)
	CUB	X ✓ Improve	$\begin{array}{c} 87.00 \pm 4.36 \\ 88.33 \pm 4.05 \\ \bigtriangleup 1.33 \end{array}$	$\begin{array}{c} 20.49 \pm 0.30 \\ 20.53 \pm 0.46 \\ \bigtriangledown 0.04 \end{array}$	$59.44 \pm 22.10 \\ 55.94 \pm 17.07 \\ \triangle 3.50$	$\begin{array}{c} 49.52 \pm 17.35 \\ 47.92 \pm 16.89 \\ \bigtriangleup 1.60 \end{array}$
Type I	Animal	X ✓ Improve	$\begin{array}{c} 81.72 \pm 2.51 \\ 82.73 \pm 1.64 \\ \bigtriangleup 1.01 \end{array}$	$\begin{array}{c} 36.87 \pm 0.41 \\ 36.87 \pm 0.36 \\ 0.00 \end{array}$	$\begin{array}{c} 82.14 \pm 27.20 \\ 71.54 \pm 16.03 \\ \bigtriangleup 10.60 \end{array}$	$\begin{array}{c} 63.94 \pm 22.74 \\ 55.50 \pm 13.13 \\ \bigtriangleup 8.44 \end{array}$
	TUAND- ROMD	X ✓ Improve	$\begin{array}{c} 84.66 \pm 0.43 \\ 85.20 \pm 0.81 \\ \bigtriangleup 0.54 \end{array}$	$\begin{array}{c} 6.88 \pm 0.00 \\ 6.88 \pm 0.00 \\ 0.00 \end{array}$	$61.46 \pm 6.09 \\ 58.24 \pm 5.05 \\ riangle 3.22$	$\begin{array}{c} 49.00 \pm 5.75 \\ 46.64 \pm 4.55 \\ \bigtriangleup 2.36 \end{array}$
	CUB	X ✓ Improve	$\begin{array}{c} 92.33 \pm 1.11 \\ 94.50 \pm 1.71 \\ \triangle 2.17 \end{array}$	$\begin{array}{c} 2.33 \pm 0.55 \\ 2.24 \pm 1.27 \\ \triangle \ 0.86 \end{array}$	$\begin{array}{c} 10.92 \pm 1.94 \\ 9.32 \pm 3.91 \\ \bigtriangleup 1.60 \end{array}$	$\begin{array}{c} 7.82 \pm 1.32 \\ 7.60 \pm 3.02 \\ \triangle \ 0.22 \end{array}$
Type II	Animal	X ✓ Improve	$\begin{array}{c} 86.75 \pm 0.33 \\ 87.61 \pm 0.50 \\ \bigtriangleup 0.86 \end{array}$	$\begin{array}{c} 8.25 \pm 3.79 \\ 4.99 \pm 0.46 \\ \bigtriangleup 3.26 \end{array}$	$\begin{array}{c} 27.62 \pm 7.42 \\ 21.26 \pm 1.31 \\ \bigtriangleup 6.36 \end{array}$	$\begin{array}{c} 18.40 \pm 7.27 \\ 13.24 \pm 0.92 \\ \bigtriangleup 5.16 \end{array}$
	TUAND- ROMD	X ✓ Improve	$86.32 \pm 0.85 \\ 88.69 \pm 0.99 \\ riangle 2.37$	$3.26 \pm 0.09 \\ 3.21 \pm 0.15 \\ riangle 0.05$	$\begin{array}{c} 43.40 \pm 2.65 \\ 38.62 \pm 5.44 \\ \bigtriangleup 4.78 \end{array}$	$33.56 \pm 2.38 \\ 31.90 \pm 4.37 \\ riangle 1.66$
	NYUD2	X ✓ Improve	$\begin{array}{c} 66.89 \pm 0.85 \\ 68.09 \pm 0.68 \\ \bigtriangleup 1.20 \end{array}$	$\begin{array}{c} 10.03 \pm 0.10 \\ 9.83 \pm 0.15 \\ \bigtriangleup 0.20 \end{array}$	$\begin{array}{c} 140.53 \pm 5.66 \\ 137.27 \pm 6.94 \\ \bigtriangleup 3.26 \end{array}$	$78.40 \pm 5.01 79.87 \pm 6.30 \triangle 1.47$
Type III	SUN- RGBD	× ✓ Improve	$\begin{array}{c} 62.11 \pm 0.31 \\ 62.78 \pm 0.32 \\ \bigtriangleup 0.67 \end{array}$	$\begin{array}{c} 13.27 \pm 0.53 \\ 13.25 \pm 0.46 \\ \bigtriangleup 0.05 \end{array}$	$\begin{array}{c} 181.00 \pm 1.20 \\ 174.90 \pm 1.50 \\ \triangle \ 6.10 \end{array}$	$\begin{array}{c} 97.87 \pm 1.48 \\ 95.00 \pm 1.00 \\ \bigtriangleup 2.87 \end{array}$

Table 2: Accuracy performance comparison for whether the model is equipped with the CML regularization term (i.e., whether λ is set to 0).

4.3.2 CML REGULARIZATION IMPROVES ROBUSTNESS

In this subsection, we evaluate the performance on the complete multimodal data, where the training/test data is divided as previous work (Zhang et al., 2019). From Tab. 2, the classification models equipped with CML regularization consistently outperform their counterpart (i.e., the original classification models) validating the rationality of CML principle. Limited by space, results on more datasets are shown in Appendix C.4.

We also find that CML regularization can improve the robustness of imperfect data, such as noise. We evaluate the models in terms of the accuracy in the test under Gaussian noise (i.e., zero mean and varying variance ϵ), and "Noise On" indicates which modality is noised (e.g., {1} indicates the first modality is noised). We report the performance on the challenging datasets (CUB and Animal) in the main text (Tab. 3) and more results are in Appendix C.3. We can find that the models equipped with CML regularization are more robust to noise, especially when the noise is much heavier.

4.3.3 PERFORMANCE UNDER DIFFERENT STRENGTHS OF CML REGULARIZATION

In this subsection, we report the accuracy under different strengths of regularization (where " $\lambda = 0$ " indicates the model is not equipped with the proposed CML regularization). We also add Gaussian noise (i.e., zero mean and varying variance ϵ) to one of the modalities, and it is clear that the model with CML regularization is more robust to the potential noise.

We show the results in Fig. 4. From Fig. 4, we can find that CML regularization can promote the accuracy on the noised data. The main reason is that the CML regularization enforces the reasonable

Dataset	Noise on	CML	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon=0.5$
CUB	{1}	× ✓ Improve	$\begin{array}{c} 84.72 \pm 3.32 \\ 85.83 \pm 2.72 \\ \bigtriangleup 1.11 \end{array}$	$\begin{array}{c} 82.22 \pm 4.53 \\ 85.00 \pm 3.50 \\ \bigtriangleup 2.78 \end{array}$	$\begin{array}{c} 79.72 \pm 4.43 \\ 84.17 \pm 4.08 \\ \bigtriangleup 4.45 \end{array}$	$\begin{array}{c} 71.17 \pm 9.14 \\ 81.11 \pm 4.37 \\ \bigtriangleup 9.94 \end{array}$
	{2}	X ✓ Improve	$\begin{array}{c} 84.44 \pm 2.75 \\ 85.83 \pm 3.40 \\ \bigtriangleup 1.39 \end{array}$	$\begin{array}{c} 83.89 \pm 3.22 \\ 85.28 \pm 2.75 \\ \bigtriangleup 1.39 \end{array}$	$\begin{array}{c} 83.61 \pm 2.83 \\ 85.28 \pm 1.97 \\ \bigtriangleup 1.67 \end{array}$	$\begin{array}{c} 83.61 \pm 3.87 \\ 85.00 \pm 1.80 \\ \bigtriangleup 1.39 \end{array}$
	{1,2}	X ✓ Improve	$\begin{array}{c} 85.00 \pm 3.12 \\ 85.83 \pm 2.72 \\ \bigtriangleup 0.83 \end{array}$	$\begin{array}{c} 82.78 \pm 3.98 \\ 85.84 \pm 3.12 \\ \bigtriangleup 3.06 \end{array}$	$\begin{array}{c} 80.00 \pm 4.46 \\ 85.83 \pm 4.25 \\ \bigtriangleup 5.83 \end{array}$	$\begin{array}{c} 72.50 \pm 11.14 \\ 81.39 \pm 6.43 \\ \bigtriangleup 8.89 \end{array}$
Animal	{1}	x ✓ Improve	$\begin{array}{c} 80.78 \pm 2.79 \\ 82.03 \pm 1.91 \\ \bigtriangleup 1.25 \end{array}$	$\begin{array}{c} 80.96 \pm 2.78 \\ 82.37 \pm 2.09 \\ \bigtriangleup 1.41 \end{array}$	$\begin{array}{c} 80.85 \pm 2.80 \\ 82.55 \pm 2.24 \\ \bigtriangleup 1.70 \end{array}$	$\begin{array}{c} 80.68 \pm 2.93 \\ 82.30 \pm 2.40 \\ \bigtriangleup 1.62 \end{array}$
	{2}	x ✓ Improve	$\begin{array}{c} 80.70 \pm 2.45 \\ 82.07 \pm 1.57 \\ \bigtriangleup 1.37 \end{array}$	$79.81 \pm 3.14 \\ 81.23 \pm 2.32 \\ \triangle 1.42$	$77.34 \pm 4.80 \\78.93 \pm 3.65 \\ \triangle 1.59$	$68.52 \pm 9.68 \\ 72.39 \pm 8.35 \\ \triangle 3.87$
	{1,2}	x ✓ Improve	$\begin{array}{c} 80.87 \pm 2.55 \\ 82.14 \pm 1.76 \\ \bigtriangleup 1.27 \end{array}$	$\begin{array}{c} 79.97 \pm 3.12 \\ 81.95 \pm 2.65 \\ \bigtriangleup 1.98 \end{array}$	$77.11 \pm 5.86 \\ 79.63 \pm 5.28 \\ \triangle 2.52$	$\begin{array}{c} 65.08 \pm 12.75 \\ 72.46 \pm 11.39 \\ \bigtriangleup 7.38 \end{array}$

Table 3: Accuracy performance comparison when some of the modalities is corrupted with Gaussian noise (i.e., zero mean with varying variance ϵ).



Figure 4: Accuracy estimation where one of the modalities is corrupted with noise.

confidence estimation. Moreover, according to Fig. 4, the proposed regularization is not sensitive to the hyperparameter λ , where quite a promising performance could be excepted with a mild regularization strength.

5 CONCLUSION

Through extensive empirical studies, we observe that the confidence estimations of current multimodal learning algorithms are typically unreliable, which tend to rely on some partial modalities. This further leads the learned model to being non-robust against the modality corruption. To be specific, model tends to be overconfident to some modalities, and ignores the evidences from other modalities even those may be useful to make decision. To solve this problem, we introduce a novel regularization technique to calibrate the confidence estimation, which forces model to learn a calibrated predictive distribution. This technique can be naturally applied to most existing multimodal learning methods without modifying their original training process and model structures. We perform comprehensive experiments to demonstrate our method's superiority in classification performance, confidence calibration and model robustness.

REFERENCES

- Moloud Abdar, Farhad Pourpanah, Sadiq Hussain, Dana Rezazadegan, Li Liu, Mohammad Ghavamzadeh, Paul Fieguth, Xiaochun Cao, Abbas Khosravi, U Rajendra Acharya, et al. A review of uncertainty quantification in deep learning: Techniques, applications and challenges. *Information Fusion*, 76:243–297, 2021.
- S. Antol, A. Agrawal, J. Lu, M. Mitchell, and D. Parikh. Vqa: Visual question answering. *Interna*tional Journal of Computer Vision, 123(1):4–31, 2015.
- AmirAli Bagher Zadeh, Paul Pu Liang, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. Multimodal language analysis in the wild: CMU-MOSEI dataset and interpretable dynamic fusion graph. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, 2018.
- Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In Proceedings of the eleventh annual conference on Computational learning theory, pp. 92–100, 1998.
- Parthajit Borah, DK Bhattacharyya, and JK Kalita. Malware dataset generation and evaluation. In 2020 IEEE 4th Conference on Information & Communication Technology (CICT), pp. 1–6. IEEE, 2020.
- Adam D Cobb and Brian Jalaian. Scaling hamiltonian monte carlo inference for bayesian neural networks with symmetric splitting. In *Uncertainty in Artificial Intelligence*, pp. 675–685. PMLR, 2021.
- Charles Corbière, Nicolas Thome, Avner Bar-Hen, Matthieu Cord, and Patrick Pérez. Addressing failure prediction by learning model confidence. In *NeurIPS*, 2019.
- John Denker and Yann LeCun. Transforming neural-net output levels to probability distributions. Advances in neural information processing systems, 3, 1990.
- Andrew Foong, David Burt, Yingzhen Li, and Richard Turner. On the expressiveness of approximate inference in bayesian neural networks. *NeurIPS*, 33:15897–15908, 2020.
- A. S. Georghiades, P. N. Belhumeur, and D. J. Kriegman. From few to many: illumination cone models for face recognition under variable lighting and pose. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 23(6):643–660, 2002.
- C. Guo, G. Pleiss, S. Yu, and K. Q. Weinberger. On calibration of modern neural networks. In *ICML*, 2017.
- Danijar Hafner, Dustin Tran, Timothy P. Lillicrap, Alex Irpan, and James Davidson. Noise contrastive priors for functional uncertainty. In UAI, 2019.
- Marton Havasi, Rodolphe Jenatton, Stanislav Fort, Jeremiah Zhe Liu, Jasper Snoek, Balaji Lakshminarayanan, Andrew M Dai, and Dustin Tran. Training independent subnetworks for robust prediction. *arXiv preprint arXiv:2010.06610*, 2020.
- Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *ICLR*, 2017.
- Hamid Reza Vaezi Joze, Amirreza Shaban, Michael L Iuzzolino, and Kazuhito Koishida. Mmtm: Multimodal transfer module for cnn fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13289–13299, 2020.
- Theofanis Karaletsos and Thang D Bui. Hierarchical gaussian process priors for bayesian neural network weights. *NeurIPS*, 33:17141–17152, 2020.
- Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? Advances in neural information processing systems, 30, 2017.
- Alireza Khodayari, Ali Ghaffari, Sina Ameli, and Jamal Flahatgar. A historical review on lateral and longitudinal control of autonomous vehicle motions. In *International Conference on Mechanical & Electrical Technology*, 2010.

- R. M. Kishi, T. H. Trojahn, and R. Goularte. Correlation based feature fusion for the temporal video scene segmentation task. *Multimedia Tools & Applications*, 78(11):15623–15646, 2019.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, volume 25, 2012.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *NeurIPS*, 2017.
- Quoc Le and Tomas Mikolov. Distributed representations of sentences and documents. In *ICML*, pp. 1188–1196. PMLR, 2014.
- Changhee Lee and Mihaela van der Schaar. A variational information bottleneck approach to multiomics data integration. In *International Conference on Artificial Intelligence and Statistics*, pp. 1513–1521. PMLR, 2021.
- David JC MacKay. Bayesian interpolation. Neural computation, 4(3):415–447, 1992.
- Pierre-Alexandre Mattei and Jes Frellsen. Miwae: Deep generative modelling and imputation of incomplete data sets. In *ICML*, pp. 4413–4423. PMLR, 2019.
- Dmitry Molchanov, Arsenii Ashukha, and Dmitry P. Vetrov. Variational dropout sparsifies deep neural networks. In ICML, 2017.
- Jishnu Mukhoti, Viveka Kulharia, Amartya Sanyal, Stuart Golodetz, Philip Torr, and Puneet Dokania. Calibrating deep neural networks using focal loss. In *NeurIPS*, 2020.
- Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help? In *NeurIPS*, 2019.
- Radford M Neal. Bayesian learning for neural networks. Springer Science & Business Media, 2012.
- S. Perkins and J. Theiler. Online feature selection using grafting. In ICML, 2003.
- Richard J Perrin, Anne M Fagan, and David M Holtzman. Multimodal techniques for diagnosis and prognosis of alzheimer's disease. *Nature*, 461(7266):916–922, 2009.
- Kefaya Qaddoum and E. L. Hines. Reliable yield prediction with regression neural networks. In *WSEAS international conference on systems theory and scientific computation*, 2012.
- Xiaojuan Qi, Renjie Liao, Jiaya Jia, Sanja Fidler, and Raquel Urtasun. 3d graph neural networks for rgbd semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5199–5208, 2017.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 567–576, 2015.
- Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. *arXiv preprint arXiv:1812.05159*, 2018.
- Joost van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation using a single deep deterministic neural network. In *ICML*, 2020.
- Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- Weiyao Wang, Du Tran, and Matt Feiszli. What makes training multi-modal classification networks hard? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12695–12705, 2020.

- Yansen Wang, Ying Shen, Zhun Liu, Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. Words can shift: Dynamically adjusting word representations using nonverbal behaviors. In *AAAI*, 2019.
- Mike Wu and Noah Goodman. Multimodal generative models for scalable weakly-supervised learning. *NeurIPS*, 31, 2018.
- Nan Wu, Stanisław Jastrzębski, Kyunghyun Cho, and Krzysztof J. Geras. Characterizing and overcoming the greedy nature of learning in multi-modal deep neural networks. In *ICML*, 2022.
- Changqing Zhang, Zongbo Han, yajie cui, Huazhu Fu, Joey Tianyi Zhou, and Qinghua Hu. Cpmnets: Cross partial multi-view networks. In *NeurIPS*, volume 32, 2019.

CONTENTS

A	How	to make ranking pairs	13
B	Ana	lysis of the training time and space complexity	14
С	Exp	eriments details	14
	C.1	Dataset details	14
	C.2	Experiment setting	14
	C.3	Robustness evaluation	16
	C.4	Additional results for robustness estimation	17
	C.5	Confidence estimation for complete inputs	17
	C.6	Confidence estimation when just penalizing the confidence difference	18
D	Algo	prithms	18
	D.1	CML for the imputation-independent model	18
	D.2	CML for the imputation-dependent model	18
E	Disc	ussion	19
	E.1	Data-imbalanced	19
	E.2	When additional modality is corrupted	20
F	CM	L being deployed in advanced multimodal models	20
G	Rela	ted work details	21

A HOW TO MAKE RANKING PAIRS



Figure 5: Illustration of generating $\mathbb S$ and $\mathbb T.$

To compute this score in practice, we initialize S as the complete modalities, and obtain T by randomly removing a modality from S. Then T is regarded as S for another confidence ranking pair and we repeat this process until there is only one modality remained in T.

B ANALYSIS OF THE TRAINING TIME AND SPACE COMPLEXITY

Ideally, CML should be computed over all possible pairs at each model update. However, it is computationally expensive, so we employ an approximation scheme following Toneva et al. (2018) for reducing the costs. For example, given samples with 4 modalities (a, b, c, d), we need to sample 3 pairs (a/ab, ab/abc, abc/abcd) to approximate CML loss, and indexes are shuffled for different epochs. So if the complexity of the traditional model is o(n), the complexity of our method will be o((k-1)n), where k indicates the number of modalities. It should be pointed out that compared models in our experiments are also equipped with sampling, and the complexity of compared methods is also o((k-1)n). We report the training time (seconds) for the same training epochs (Platform: RTX 3090×8, CUDA Version: 11.2). It is observed that the original model and model equipped with CML have the same level of computational complexity.

Table 4: Training time (Platform: RTX 3090×8).

Method	CML	TUANDROMD	YaleB	Handwritten	CUB	Animal
Type I	×	$245.3 \\ 297.6$	$1574.6 \\ 1210.2$	$\begin{array}{c} 141.5\\ 191.2 \end{array}$	$\begin{array}{c} 351.6\\ 348.5\end{array}$	$1582.7 \\ 1641.3$
Type II	×	$1447.7 \\ 1489.1$	703.3 662.9	$233.2 \\ 210.8$	$565.2 \\ 781.7$	717.8 720.3

C EXPERIMENTS DETAILS

C.1 DATASET DETAILS

We evaluate the proposed method on diverse datasets, including data with multiple modalities and multiple types of features. \circ **YaleB**: Similar to previous work Georghiades et al. (2002), we also use a subset of this face image dataset, which contains 650 facial images, 10 classes and 3 different types of features. \circ **Handwritten** Perkins & Theiler (2003): This is a database of handwritten digits which contains 2,000 images, 10 classes, 6 types of features. \circ **CUB** Wah et al. (2011): Following CPM-Nets Zhang et al. (2019), we use a subset of this dataset, which contains first 10 classes of original dataset and 2 modalities (deep visual feature and text feature) are obtained by GoogleNet and doc2vec Le & Mikolov (2014). \circ **Animal**: This dataset contains 10, 158 images, 50 classes, and 2 types of features (deep visual feature from DECAF Krizhevsky et al. (2012) and VGG19 Simonyan & Zisserman (2014)). \circ **TUANDROMD** Borah et al. (2020): The dataset contains 4, 465 instances, 2 classes and 2 types of modalities.

C.2 EXPERIMENT SETTING

Type-I: For CPM-Nets and the first five datasets(i.e.,YaleB, Handwritten, CUB and Animal), we follow the author's implementation Zhang et al. (2019): the dimensionality of latent representation is 150. Parameter lambda for cub/animal/hand-written/yaleB/tuandromd is set as 5/45/45/10/5. The dimensionalities of input, hidden layers are 128 and 300. We use Adam optimizer to train all CPM-Nets models with the learning rate of 10^{-2} and no additional regularization term. For Tuandromd dataset, we tune the dimensionality of latent representation to 512. The dimensionalities of input and hidden layers are both 512. We use Adam optimizer to train CPM-Net with L2-regularization term. **Type-II**: For MIWAE, we train the encoder, decoder and classifier respectively. The number of hidden units of them is all 128. Parameter lambda for cub/animal/hand-written/yaleB/tuandromd are set as 15/25/10/35/75 for best performance. The dimensionalities of the latent space are 64. We use Adam optimizer to train the encoder and decoder with a learning rate of 10^{-2} . Then we train the encoder, decoder and classifier rate of 10^{-2} . Then we train the encoder, decoder and classifier altogether for another with a learning rate of 10^{-3} . As same as

Dataset	Noise on	CML	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.4$	$\epsilon = 0.5$
	{1}	× ✓	97.43 \pm 1.58 98.46 \pm 1.09	$\begin{array}{c} 96.92 \pm 1.88 \\ \textbf{98.20} \pm \textbf{1.31} \end{array}$	$\begin{array}{c} 96.41 \pm 2.20 \\ \textbf{96.15 \pm 1.88} \end{array}$	$94.10 \pm 1.31 \\ \textbf{94.62 \pm 1.88}$	$\begin{array}{c} 92.82 \pm 1.31 \\ \textbf{93.59 \pm 1.30} \end{array}$
	{2}	× ✓	95.13 \pm 0.72 96.92 \pm 1.26	$\begin{array}{c} 94.10 \pm 1.31 \\ \textbf{95.90} \pm \textbf{2.02} \end{array}$	$\begin{array}{c} 92.57 \pm 0.73 \\ \textbf{94.61 \pm 2.88} \end{array}$	$\begin{array}{c} 92.05 \pm 1.45 \\ \textbf{93.33} \pm \textbf{2.54} \end{array}$	$\begin{array}{c} 91.54 \pm 1.66 \\ \textbf{93.08} \pm \textbf{3.14} \end{array}$
	{3}	× ✓	$\begin{array}{c} 94.87 \pm 0.96 \\ \textbf{96.92 \pm 1.88} \end{array}$	$\begin{array}{c} 94.87 \pm 0.96 \\ \textbf{97.18} \pm \textbf{1.92} \end{array}$	$\begin{array}{c} 94.10 \pm 0.96 \\ \textbf{96.15 \pm 1.88} \end{array}$	$\begin{array}{c} 92.82 \pm 1.81 \\ \textbf{94.87} \pm \textbf{2.54} \end{array}$	$\begin{array}{c} 92.05 \pm 1.31 \\ \textbf{94.36} \pm \textbf{2.02} \end{array}$
YaleB	{1,2}	× ✓	96.67 ± 2.61 97.69 ± 0.63	95.13 ± 3.46 95.39 ± 2.26	$\begin{array}{c} 91.28 \pm 2.83 \\ \textbf{92.56} \pm \textbf{2.02} \end{array}$	$\begin{array}{c} 88.72 \pm 3.10 \\ \textbf{89.72} \pm \textbf{2.21} \end{array}$	86.41 ± 3.10 86.66 \pm 1.81
	{1,3}	× ✓	97.43 \pm 0.96 98.46 \pm 1.09	97.69 ± 1.66 98.46 ± 1.26	97.43 ± 1.81 98.46 \pm 1.66	97.18 ± 2.20 96.92 ± 1.88	96.15 ± 2.26 96.67 ± 2.20
	{2,3}	× ✓	$\begin{array}{ c c c c c } 94.62 \pm 1.08 \\ \textbf{96.41} \pm \textbf{1.81} \end{array}$	$\begin{array}{c} 93.85 \pm 1.25 \\ \textbf{95.64 \pm 1.92} \end{array}$	90.26 ± 2.54 93.84 ± 3.32	87.95 ± 2.83 91.28 \pm 3.10	86.67 ± 2.38 89.49 \pm 3.16
	{1, 2, 3}	× ✓	96.15 \pm 1.88 97.43 \pm 1.81	96.41 ± 3.16 97.43 \pm 1.92	93.85 ± 4.40 93.85 ± 4.40	87.69 ± 8.21 87.69 \pm 7.61	$84.10 \pm 10.32 \\ \textbf{82.56} \pm \textbf{9.26}$
	{1}	× ✓	97.18 ± 1.92 98.46 ± 1.26	95.38 ± 1.25 95.90 ± 1.92	93.34 ± 1.31 93.85 ± 1.88	92.57 ± 1.58 93.08 ± 1.66	91.28 ± 1.31 92.31 ± 0.63
	{2}	× ✓	88.46 ± 1.66 90.77 ± 3.33	87.18 ± 1.31 90.26 \pm 3.57	86.92 ± 1.09 89.75 \pm 3.85	86.92 ± 1.09 89.75 \pm 3.84	$\begin{array}{c} 86.92 \pm 1.09 \\ \textbf{89.75} \pm \textbf{3.84} \end{array}$
	{3}	× ✓	85.90 ± 1.92 88.97 ± 2.54	85.13 ± 1.81 88.21 \pm 2.61	84.87 ± 1.45 87.69 ± 2.74	84.62 ± 1.66 87.69 ± 3.32	$\begin{array}{c} 84.62 \pm 1.66 \\ \textbf{87.44} \pm \textbf{3.10} \end{array}$
Hand- written	{1,2}	× ✓	88.97 ± 3.68 88.97 ± 4.04	83.08 ± 3.50 83.59 ± 2.97	$\begin{array}{c} 78.97 \pm 1.92 \\ \textbf{80.51} \pm \textbf{3.46} \end{array}$	77.69 ± 2.74 77.18 \pm 4.28	75.90 ± 3.57 74.10 \pm 3.84
	{1,3}	× ✓	91.54 \pm 1.09 93.59 \pm 2.38	91.28 ± 3.16 91.79 ± 3.68	88.97 ± 5.41 88.97 \pm 4.04	87.43 ± 5.83 86.93 \pm 4.99	85.64 ± 6.42 85.39 ± 4.91
	{2,3}	X V	$ \begin{array}{c} 63.59 \pm 8.00 \\ \mathbf{64.36 \pm 7.49} \end{array} $	$59.74 \pm 7.00 \\ 58.46 \pm 6.37$	57.69 ± 5.99 56.67 ± 6.10	$56.67 \pm 5.94 \\ 55.64 \pm 6.04$	55.90 ± 5.49 54.87 ± 6.29
	{1, 2, 3}	× ✓	$\begin{array}{c c} 54.87 \pm 10.68 \\ \textbf{57.18} \pm \textbf{11.41} \end{array}$	37.95 ± 6.92 35.64 ± 4.80	$29.48 \pm 4.76 \\ 26.67 \pm 2.54$	$24.36 \pm 4.04 \\ 22.82 \pm 2.54$	$22.31 \pm 4.12 \\20.77 \pm 1.09$
	{1}	× ✓	84.72 ± 3.32 85.83 ± 2.72	$\begin{array}{c} 82.22 \pm 4.53 \\ \textbf{85.00 \pm 3.50} \end{array}$	$\begin{array}{c} 79.72 \pm 4.43 \\ \textbf{84.17 \pm 4.08} \end{array}$	$\begin{array}{c} 76.39 \pm 6.85 \\ \textbf{83.06 \pm 3.99} \end{array}$	$\begin{array}{c} 71.17 \pm 9.14 \\ \textbf{81.11} \pm \textbf{4.37} \end{array}$
CUB	{2}	× ✓	$\begin{array}{c c} 84.44 \pm 2.75 \\ \textbf{85.83} \pm \textbf{3.40} \end{array}$	$\begin{array}{c} 83.89 \pm 3.22 \\ \textbf{85.28} \pm \textbf{2.75} \end{array}$	$\begin{array}{c} 83.61 \pm 2.83 \\ \textbf{85.28} \pm \textbf{1.97} \end{array}$	83.89 ± 3.49 85.28 ± 1.97	$\begin{array}{c} 83.61 \pm 3.87 \\ \textbf{85.00 \pm 1.80} \end{array}$
	{1,2}	× ✓	85.00 ± 3.12 85.83 ± 2.72	$\begin{array}{c} 82.78 \pm 3.98 \\ \textbf{85.84 \pm 3.12} \end{array}$	80.00 ± 4.46 85.83 \pm 4.25	$\begin{array}{c} 76.67 \pm 7.48 \\ \textbf{84.44} \pm \textbf{4.38} \end{array}$	$\begin{array}{c} 72.50 \pm 11.14 \\ \textbf{81.39} \pm \textbf{6.43} \end{array}$
	{1}	× ✓	$ \begin{array}{c c} 80.78 \pm 2.79 \\ $	$\begin{array}{c} 80.96 \pm 2.78 \\ \textbf{82.37} \pm \textbf{2.09} \end{array}$	$\begin{array}{c} 80.85 \pm 2.80 \\ \textbf{82.55} \pm \textbf{2.24} \end{array}$	$\begin{array}{c} 80.81 \pm 2.88 \\ \textbf{82.42} \pm \textbf{2.22} \end{array}$	$\begin{array}{c} 80.68 \pm 2.93 \\ \textbf{82.30} \pm \textbf{2.40} \end{array}$
Animal	{2}	× ✓	80.70 \pm 2.45 82.07 \pm 1.57	$\begin{array}{c} 79.81 \pm 3.14 \\ \textbf{81.23} \pm \textbf{2.32} \end{array}$	$77.34 \pm 4.80 \\ \textbf{78.93} \pm \textbf{3.65}$	$\begin{array}{c} 72.89 \pm 7.46 \\ \textbf{75.81} \pm \textbf{6.30} \end{array}$	$68.52 \pm 9.68 \\ \textbf{72.39} \pm \textbf{8.35}$
	{1,2}	× ✓	80.87 ± 2.55 82.14 ± 1.76	$\begin{array}{c} 79.97 \pm 3.12 \\ \textbf{81.95} \pm \textbf{2.65} \end{array}$	77.11 ± 5.86 79.63 ± 5.28	$\begin{array}{c} 72.23 \pm 9.04 \\ \textbf{76.63 \pm 7.73} \end{array}$	$\begin{array}{c} 65.08 \pm 12.75 \\ \textbf{72.46} \pm \textbf{11.39} \end{array}$
	{1}	× ✓	$\begin{array}{c} 84.77 \pm 0.55 \\ \textbf{86.50} \pm \textbf{0.59} \end{array}$	$\begin{array}{c} 80.47 \pm 0.99 \\ \textbf{82.46} \pm \textbf{0.77} \end{array}$	$\begin{array}{c} 76.53 \pm 1.11 \\ \textbf{78.30} \pm \textbf{1.18} \end{array}$	$\begin{array}{c} 72.65 \pm 0.76 \\ \textbf{74.92 \pm 1.39} \end{array}$	$\begin{array}{c} 70.17 \pm 0.66 \\ \textbf{72.45} \pm \textbf{1.33} \end{array}$
TUAND- ROMD	{2}	× ✓	86.56 ± 0.27 88.87 ± 0.22	$\begin{array}{c} 85.71 \pm 0.48 \\ \textbf{88.74 \pm 0.28} \end{array}$	$84.14 \pm 0.58 \\ \textbf{88.58} \pm \textbf{0.63}$	$\begin{array}{c} 82.35 \pm 0.86 \\ \textbf{88.15} \pm \textbf{0.65} \end{array}$	80.85 ± 1.05 87.93 \pm 0.67
	{1,2}	×	84.88 ± 1.19 87.41 ± 3.40	80.72 ± 1.02 82.78 ± 1.14	76.60 ± 0.75 79.28 \pm 1.00	73.15 ± 1.10 76.30 \pm 1.11	70.35 ± 1.25 73.82 \pm 1.35

Table 5: Accuracy performance comparison when some of the modalities is blurred (Type I).

prior work Corbière et al. (2019), we evaluated the performance according to Accuracy (%), NLL (10^{-1}), AURC (10^{-3}), and E-AURC (10^{-3}). For both types above, we set 0 as the default value of τ to ensure the model meets CML strictly.

Dataset	Noise Noise on	CML	$\epsilon = 0.5$	$\epsilon = 1.0$	$\epsilon = 1.5$	$\epsilon = 2.0$	$\epsilon = 2.5$
	{1}	×	95.90 ± 2.54 97.43 ± 1.31	94.87 ± 3.22 96.15 ± 2.51	93.85 ± 2.88 95.13 ± 2.97	93.59 ± 3.16 94.36 ± 2.97	93.59 ± 3.16 93.85 ± 3.46
	{2}	×	96.15 \pm 2.26 97.69 \pm 1.26	93.33 ± 3.22 96.67 ± 1.58	91.03 ± 2.62 94.10 ± 2.20	90.26 ± 2.02 92.82 ± 2.83	89.23 ± 2.18 92.05 ± 2.02
	{3}	×	98.72 \pm 0.36 98.72 \pm 0.73	96.92 ± 1.26 97.69 ± 1.09	96.15 ± 0.63 97.43 \pm 0.96	96.15 ± 0.63 97.18 \pm 1.31	95.90 ± 0.96 96.67 \pm 1.58
YaleB	{1, 2}	× ✓	95.64 \pm 2.83 96.66 \pm 1.31	91.02 ± 3.46 93.59 ± 2.38	88.46 ± 4.53 90.51 \pm 2.97	87.18 ± 3.46 86.67 \pm 3.46	85.90 ± 4.09 84.62 \pm 3.26
	{1, 3}	X V	98.46 ± 0.63 98.20 ± 0.73	98.46 ± 1.66 97.95 ± 1.92	97.69 ± 1.66 97.69 ± 1.66	97.43 ± 1.45 98.20 ± 1.58	97.18 ± 1.31 97.69 ± 1.66
	{2, 3}	× ✓	97.43 \pm 0.36 98.72 \pm 0.36	95.89 ± 0.36 97.69 ± 1.09	95.38 ± 0.62 96.66 ± 0.73	$\begin{array}{c} 94.62 \pm 0.62 \\ \textbf{95.38} \pm \textbf{0.62} \end{array}$	$\begin{array}{c} 92.82 \pm 0.73 \\ \textbf{94.61} \pm \textbf{1.66} \end{array}$
	{1, 2, 3}	× ✓	97.69 ± 0.63 98.46 ± 0.63	95.64 ± 0.36 97.18 ± 1.31	93.08 ± 1.09 95.64 ± 0.96	89.23 ± 1.66 92.56 ± 2.54	82.31 ± 1.26 88.46 \pm 2.27
	{1}	× ✓	$98.42 \pm 0.51 \\ \textbf{99.50} \pm \textbf{0.41}$	$\begin{array}{c} 98.25 \pm 0.35 \\ \textbf{99.50} \pm \textbf{0.41} \end{array}$	97.92 ± 0.12 99.50 ± 0.41	$\begin{array}{c} 97.92 \pm 0.12 \\ \textbf{99.50} \pm \textbf{0.41} \end{array}$	97.58 ± 0.12 99.50 ± 0.41
	{2}	× ✓	98.17 \pm 1.03 98.83 \pm 0.24	97.75 ± 0.54 98.50 ± 0.41	97.33 ± 0.47 98.67 ± 0.47	97.00 ± 0.41 98.67 ± 0.47	96.92 ± 0.42 98.67 ± 0.47
Hand- written	{1,2}	× ✓	97.67 \pm 0.47 99.00 \pm 0.00	97.25 ± 0.54 98.83 ± 0.24	96.58 ± 0.51 98.83 ± 0.24	95.92 ± 0.59 98.83 ± 0.24	95.67 ± 0.94 98.83 ± 0.24
WINCOM	{1,3}	× ✓	$\begin{array}{c} 98.08 \pm 0.12 \\ \textbf{99.50} \pm \textbf{0.00} \end{array}$	97.00 ± 1.22 99.17 ± 0.47	96.33 ± 1.55 98.00 ± 0.41	95.33 ± 1.55 97.67 ± 0.24	$\begin{array}{c} 95.08 \pm 1.59 \\ \textbf{95.78} \pm \textbf{1.04} \end{array}$
	{2, 3}	× ✓	98.17 \pm 0.24 99.00 \pm 0.00	96.83 ± 0.47 98.67 ± 0.47	95.67 ± 0.85 97.67 ± 0.62	$\begin{array}{c} 94.75 \pm 0.94 \\ \textbf{97.33} \pm \textbf{1.03} \end{array}$	$\begin{array}{c} 94.17 \pm 0.85 \\ \textbf{97.33} \pm \textbf{1.03} \end{array}$
	{1, 2, 3}	X V	96.50 \pm 1.08 98.50 \pm 0.71	$\begin{array}{c} 93.58 \pm 1.74 \\ \textbf{97.17} \pm \textbf{1.18} \end{array}$	$\begin{array}{c} 90.67 \pm 3.09 \\ \textbf{95.50} \pm \textbf{1.22} \end{array}$	$\begin{array}{c} 88.75 \pm 3.54 \\ \textbf{93.67} \pm \textbf{0.85} \end{array}$	87.58 ± 3.36 92.50 \pm 0.82
	{1}	X V	91.11 \pm 1.04 93.33 \pm 1.80	86.94 ± 2.83 90.83 \pm 2.45	83.61 ± 3.93 87.50 ± 3.60	80.83 ± 4.14 85.56 \pm 4.38	$\begin{array}{c} 79.17 \pm 3.79 \\ \textbf{81.11} \pm \textbf{4.53} \end{array}$
CUB	{2}	× ✓	91.11 \pm 0.40 93.61 \pm 1.04	$\begin{array}{c} 91.95 \pm 0.39 \\ \textbf{92.78} \pm \textbf{1.04} \end{array}$	91.11 ± 0.40 92.50 ± 1.80	89.72 ± 0.39 91.67 ± 2.96	$88.61 \pm 0.79 \\ \textbf{91.39} \pm \textbf{3.22}$
	{1, 2}	X V	92.78 \pm 1.97 94.72 \pm 2.19	88.61 ± 1.42 92.22 \pm 3.75	85.83 ± 1.80 90.00 ± 4.46	$\begin{array}{c} 79.72 \pm 2.83 \\ \textbf{86.11 \pm 4.10} \end{array}$	$74.17 \pm 4.46 \\ \textbf{79.17} \pm \textbf{4.91}$
	{1}	× ✓	$\begin{array}{c} 86.61 \pm 0.20 \\ \textbf{87.20} \pm \textbf{0.18} \end{array}$	85.81 ± 0.36 87.01 \pm 0.18	$\begin{array}{c} 84.82 \pm 1.02 \\ \textbf{86.60} \pm \textbf{0.20} \end{array}$	$\begin{array}{c} 83.77 \pm 1.29 \\ \textbf{86.03} \pm \textbf{0.04} \end{array}$	$\begin{array}{c} 82.16 \pm 2.32 \\ \textbf{85.42 \pm 0.29} \end{array}$
Animal	{2}	× ✓	$\begin{array}{c} 86.33 \pm 0.54 \\ \textbf{87.04 \pm 0.08} \end{array}$	85.62 ± 0.61 86.64 \pm 0.26	84.84 ± 0.95 85.95 \pm 0.42	83.04 ± 1.24 84.78 \pm 0.17	81.34 ± 1.73 82.71 \pm 0.24
	{1, 2}	X V	$\begin{array}{c} 86.01 \pm 0.17 \\ \textbf{87.04} \pm \textbf{0.42} \end{array}$	$\begin{array}{c} 84.80 \pm 0.81 \\ \textbf{86.50} \pm \textbf{0.15} \end{array}$	83.17 ± 1.65 85.38 ± 0.34	80.92 ± 2.77 83.84 \pm 0.65	$77.42 \pm 4.14 \\ \textbf{81.67} \pm \textbf{0.75}$
	{1}	× ✓	$81.14 \pm 0.70 \\ 81.99 \pm 1.99$	$78.21 \pm 0.92 \\ \textbf{78.79} \pm \textbf{2.42}$	$\begin{array}{c} 75.39 \pm 1.09 \\ \textbf{76.37} \pm \textbf{2.57} \end{array}$	$73.21 \pm 1.46 \\ \textbf{74.36} \pm \textbf{2.63}$	$71.71 \pm 1.26 \\ \textbf{73.19} \pm \textbf{2.60}$
TUAND- ROMD	{2}	× ✓	$84.19 \pm 0.82 \\ 84.88 \pm 1.62$	84.43 ± 0.48 84.73 ± 1.89	84.46 ± 0.35 84.84 ± 1.76	84.32 ± 0.45 84.39 ± 0.89	$84.21 \pm 0.44 \\ \textbf{84.97} \pm \textbf{1.52}$
	{1,2}	X V	83.56 ± 1.23 83.99 ± 1.87	80.85 ± 1.30 81.48 ± 2.30	77.85 ± 1.53 78.50 \pm 2.30	75.90 ± 2.07 76.73 \pm 2.19	74.08 ± 2.22 75.23 \pm 2.20

Table 6: Accuracy performance comparison when some of the modalities is blurred (Type II).

C.3 ROBUSTNESS EVALUATION

We evaluate models in terms of accuracy under Gaussian noise (i.e., zero mean and varying variance ϵ), and "Noise On" indicates which modality is noised (e.g., {1} indicates the first modality is noised). In addition to the performance on the challenging datasets (CUB and Animal) in the main text (Tab. 3), we show more other results (Tab. 5 6). It is clear that the models equipped with CML are more robust to noise, especially when the noise is much heavier.

C.4 Additional results for robustness estimation

Method	Dataset	CML	Accuracy (†)	NLL (↓)	AURC (↓)	E-AURC (↓)
Type I	YaleB	× ✓ Improve	$\begin{array}{c} 95.84 \pm 0.78 \\ 97.69 \pm 1.09 \\ \bigtriangleup 1.85 \end{array}$	$\begin{array}{c} 21.98 \pm 0.05 \\ 21.98 \pm 0.05 \\ 0.00 \end{array}$	$\begin{array}{c} 3.00 \pm 1.38 \\ 1.46 \pm 1.51 \\ riangle 1.54 \end{array}$	$2.08 \pm 1.37 \\ 1.12 \pm 1.32 \\ riangle 0.96$
	Hand- written	x ✓ Improve	$\begin{array}{c} 89.00 \pm 3.64 \\ 93.60 \pm 0.60 \\ \bigtriangleup 4.60 \end{array}$	$\begin{array}{c} 20.30 \pm 0.25 \\ 20.06 \pm 0.11 \\ \bigtriangleup 0.14 \end{array}$	$\begin{array}{c} 35.83 \pm 20.43 \\ 11.00 \pm 6.17 \\ \bigtriangleup 14.83 \end{array}$	$\begin{array}{c} 28.80 \pm 15.49 \\ 8.90 \pm 5.80 \\ \bigtriangleup 19.90 \end{array}$
Type II	YaleB	x ✓ Improve	$\begin{array}{c} 95.69 \pm 2.10 \\ 97.84 \pm 0.58 \\ \bigtriangleup 2.15 \end{array}$	$\begin{array}{c} 1.80 \pm 0.71 \\ 1.11 \pm 0.49 \\ \bigtriangleup 0.69 \end{array}$	5.50 ± 2.86 5.02 ± 6.39 $\triangle 0.48$	$\begin{array}{c} 4.32 \pm 2.32 \\ 4.76 \pm 6.26 \\ \bigtriangledown 0.44 \end{array}$
	Hand- written	× ✓ Improve	$\begin{array}{c} 98.40 \pm 0.64 \\ 99.05 \pm 0.19 \\ \bigtriangleup 0.65 \end{array}$	$\begin{array}{c} 0.49 \pm 0.12 \\ 0.50 \pm 0.10 \\ 0.00 \end{array}$	$0.32 \pm 0.16 \\ 0.18 \pm 0.07 \\ riangle 0.14$	$\begin{array}{c} 0.16 \pm 0.12 \\ 0.14 \pm 0.08 \\ \triangle \ 0.02 \end{array}$

Table 7: Accuracy performance comparison for whether the model is equipped with the cma regularization term on additional dataset (i.e., whether λ is set to 0).

Limited by space, we show the performance of model equipped with CML on YaleB and Handwritten. From Table 7, the classification models equipped with CML consistently outperforms their counterpart validating the rationality of CML principle.



C.5 CONFIDENCE ESTIMATION FOR COMPLETE INPUTS

Figure 6: Confidence estimation on complete inputs. We estimate the confidence on complete inputs (top) and the confidence when one modality is removed (bottom). We can find CML regularization keeps the confidence estimation on complete input but alleviate the over-confidence when one modality is removed, which indicates the proposed method calibrates the multimodal model by rethinking the relationship between the modalities.

We show the confidence estimation for complete inputs, as shown in Fig. 6, we can find that the confidence estimation of original model and CML model are very similar. To prevent the model from being over-confident when model predicts a wrong prediction, the regularization will not be added when prediction of complete input is wrong. From the bottom figures, we can find CML regularization alleviates the problem that model increases the confidence when one modality is removed.

Proof of Lemma 3.2: if we have $\operatorname{VRR}_{CML} < \operatorname{VRR}_{Ori}$, then we have $\mathbb{E}\left(\operatorname{Conf}_{CML}(\boldsymbol{x}^{(\mathbb{T})})\right) - \mathbb{E}\left(\operatorname{Conf}_{CML}(\boldsymbol{x}^{(\mathbb{S})})\right) \leq \mathbb{E}\left(\operatorname{Conf}_{Ori}(\boldsymbol{x}^{(\mathbb{T})})\right) - \mathbb{E}\left(\operatorname{Conf}_{Ori}(\boldsymbol{x}^{(\mathbb{S})})\right)$, then we have:

$$\mathbb{E}\left(\operatorname{Conf}_{CML}(\boldsymbol{x}^{(\mathbb{T})})\right) \leq \mathbb{E}\left(\operatorname{Conf}_{Ori}(\boldsymbol{x}^{(\mathbb{T})})\right),$$
subject to: $\mathbb{E}\left(\operatorname{Conf}_{CML}(\boldsymbol{x}^{(\mathbb{T})})\right) = \mathbb{E}\left(\operatorname{Conf}_{Ori}(\boldsymbol{x}^{(\mathbb{T})})\right)$
(7)

During the train stage, we evaluate the confidence difference between the $\mathbb{E}\left(\operatorname{Conf}_{CML}(\mathbf{x}^{(\mathbb{T})})\right)$ and $\mathbb{E}\left(\operatorname{Conf}_{Ori}(\mathbf{x}^{(\mathbb{T})})\right)$, i.e., $\mathbb{E}\left(\left|\operatorname{Conf}_{CML}(\mathbf{x}^{(\mathbb{T})}) - \operatorname{Conf}_{Ori}(\mathbf{x}^{(\mathbb{T})})\right|\right)$. We find the confidence difference between the $\mathbb{E}\left(\operatorname{Conf}_{CML}(\mathbf{x}^{(\mathbb{T})})\right)$ and $\mathbb{E}\left(\operatorname{Conf}_{Ori}(\mathbf{x}^{(\mathbb{T})})\right)$ is very small (less than 0.1%), which implies that the confidence estimation on complete inputs are very close.

C.6 CONFIDENCE ESTIMATION WHEN JUST PENALIZING THE CONFIDENCE DIFFERENCE



Figure 7: Confidence estimation when penalizing the confidence difference.

Forcing the confidence for $\mathbf{x}^{(\mathbb{T})}$ to be smaller than the confidence for $\mathbf{x}^{(\mathbb{S})}$ strictly (Eq. 3) will lead to a very small confidence for $\mathbf{x}^{(\mathbb{T})}$ and will make the model estimate an extremely small confidence for each modality, which contradicts the fact that the model sometimes can still make correct predictions confidently when one modality is removed. A flexible ranking regularization makes it more suitable for real data.

D ALGORITHMS

In addition to the general algorithm shown in the main text, we show the specific algorithms corresponding to different types of algorithms and add more comments for better understanding.

D.1 CML FOR THE IMPUTATION-INDEPENDENT MODEL

D.2 CML FOR THE IMPUTATION-DEPENDENT MODEL

For imputation-dependent method, we use MIWAE to train the reconstruction model first, then we use the reconstructed modalities to train the classifier.

For reconstruction-based method, the missing modalities need to be reconstructed first, so the process can be divided into two stages.

Algorithm 2: CML for the imputation-independent model

1 Given dataset $\mathcal{D} = \left\{ \{x_i^m\}_{m=1}^M, y_i\}_{i=1}^N$, classifier f, and classification loss function \mathcal{L}^{cl} , Coefficient λ of CML, epochs for training the classifier *epoch*

2 for $e = 1, \dots, epoch$ do 3 $\mid \mathbb{S} \leftarrow \mathbb{M}$

4 Make the prediction via input \mathbb{S}

 $\begin{array}{c|c} \mathcal{L}^{\mathrm{cl}} \leftarrow \mathcal{L}^{\mathrm{cl}}(\mathbf{x}^{(\mathbb{S})}) \\ \mathcal{L}^{\mathrm{CML}} \leftarrow 0 \text{ for } m = M - 1, \dots, 1 \text{ do} \end{array}$

Randomly erase a modality of \mathbb{S} and set it as \mathbb{T}

Make the prediction via input \mathbb{T}

9 $\mathcal{L}^{cl} \leftarrow \mathcal{L}^{cl} + \mathcal{L}^{cl}(\mathbf{x}^{(\mathbb{T})})$

 $\mathcal{L}^{\text{CML}} \leftarrow \mathcal{L}^{\text{CML}} + \max\left(0, \text{Conf}(\mathbf{x}^{(\mathbb{T})}) - \text{Conf}(\mathbf{x}^{(\mathbb{S})}) - \tau\right)$

11 end

5

6

7

8

10

12 $\mathcal{L} = \frac{1}{M} \mathcal{L}^{cl} + \lambda \mathcal{L}^{CML}$

13 Update the parameters of the classification model with \mathcal{L}

14 end

15 **Return** the classifier f_{cl}

Algorithm 3: CML for the imputation-dependent model

Given dataset $\mathcal{D} = \{\{x_i^m\}_{m=1}^M, y_i\}_{i=1}^N$, reconstruction network f_{re} and classifier f_{cl} , reconstruction loss function \mathcal{L}^{re} , Coefficient λ of CML, epochs for training the reconstruction net $epoch_{re}$ and classifier $epoch_{cl}$ 2 for $e_1 = 1, ..., epoch_{re}$ do Reconstruct the modalities via reconstruction model 3 Compute the reconstruction loss by \mathcal{L}^{re} 4 Update the parameters of the reconstruction model 5 6 end 7 for $e_2 = 1, \ldots, epoch_{cl}$ do $\mathbb{S} \gets \mathbb{M}$ 8 $\mathcal{L}^{\text{CE}} \leftarrow \mathcal{L}^{\text{CE}}(\mathbf{x}^{(\mathbb{S})})$ $\mathcal{L}^{\text{CML}} \leftarrow 0$ 10 for m = M - 1, ..., 1 do 11 Randomly erase a modality of \mathbb{S} and set it as \mathbb{T} 12 Reconstruct the erased modalities via reconstruction model and add them to $\mathbf{x}^{(\mathbb{T})}$ 13 Compute the classification loss $\mathcal{L}^{CE}(\mathbf{x}^{(\mathbb{T})})$ with Cross-Entropy loss function 14 $\mathcal{L}^{\text{CE}} \leftarrow \mathcal{L}^{\text{CE}} + \mathcal{L}^{\text{CE}}(\mathbf{x}^{(\mathbb{T})})$ 15 $\mathcal{L}^{\text{CML}} \leftarrow \mathcal{L}^{\text{CML}} + \max\left(0, \text{Conf}(\mathbf{x}^{(\mathbb{T})}) - \text{Conf}(\mathbf{x}^{(\mathbb{S})}) - \tau\right)$ 16 end 17

18 $\mathcal{L} = \frac{1}{M} \mathcal{L}^{CE} + \lambda \mathcal{L}^{CML}$

19 Update the parameters of the classification model with \mathcal{L}

```
20 end
```

21 **Return** the reconstruction model f_{re} and classifier f_{cl}

E DISCUSSION

E.1 DATA-IMBALANCED

• Why the CML can still work when the training data is data-imbalanced (e.g., long-tailed)?

CML can improve performance when the data for the training model is data-imbalanced since it increases the confidence of the minority classes. For a trustworthy model, the model should treat the majority and minority classes equally during the test. CML requires the model to make predictions

fairly regardless of whether the majority and minority classes of the samples belong. On the contrary, the original model tends to predict lower confidence for the minority classes than the majority classes. And the improvements on the data-imbalanced dataset Animal (data distribution is shown in Fig. 8) validate the effectiveness.



Figure 8: Illustration of data distribution of Animal dataset (the number of samples for every classes).

Animal is a data-imbalanced real-world dataset, the improvement shows CML can also deal with applications that suffer from data-imbalanced. The original model tends to predict lower confidence for the minority classes than the majority classes, which is unfair to minority classes. CML requires the model to make predictions fairly regardless of whether the majority and minority classes of the samples belong.

E.2 WHEN ADDITIONAL MODALITY IS CORRUPTED

• Why the confidence should not decrease even when the additional modality is corrupted?

Intuitively, the predictive confidence should decrease when some features are corrupted:

$$\operatorname{Conf}(\mathbf{x}) > \operatorname{Conf}(\mathbf{x}_{\epsilon}),$$
 (8)

where \mathbf{x}_{ϵ} indicates the input \mathbf{x} is corrupted by the perturbation ϵ . It seems CML is no longer work, and the confidence relationship between the $\mathbf{x}^{(\mathbb{T})}$ and $\mathbf{x}^{(\mathbb{S})}$ should be changed:

$$\operatorname{Conf}(\mathbf{x}^{(\mathbb{T})}) > \operatorname{Conf}(\mathbf{x}^{(\mathbb{S})}).$$
(9)

Mathematically, however, Eq. 9 is not the multi-modality representation of Eq. 8, and the Eq. 10 should hold:

$$\operatorname{Conf}(\mathbf{x}^{(\mathbb{S})}) = \operatorname{Conf}(\mathbf{x}^{(\mathbb{T})} \cup \mathbf{x}^{(\mathcal{L}_{\mathbb{S}}\mathbb{T})}) > \operatorname{Conf}(\mathbf{x}^{(\mathbb{T})} \cup \mathbf{x}^{(\mathcal{L}_{\mathbb{S}}\mathbb{T})}_{\epsilon})$$
(10)

Theoretically, the confidence decrease is in terms of $\mathbf{x}^{(\mathbb{S})}$ rather than $\mathbf{x}^{(\mathbb{T})}$ when the additional modalities $\mathbf{x}^{(\mathbb{G}_{\mathbb{S}}\mathbb{T})}$ are corrupted.

F CML BEING DEPLOYED IN ADVANCED MULTIMODAL MODELS

MMTM is a SOTA method in multimodal learning which is selected as a representative method by Wu et al. (2022) and originally proposed by Joze et al. (2020). NYU Depth V2 and SUN RGB-D are two widely used multi-modality datasets for RGB-D scene recognition. \circ NYUD2: Following previous work Georghiades et al. (2002), we use a reorganized version of this dataset, which contains 1449 samples, 10 scene classes. \circ SUN RGB-D Perkins & Theiler (2003): This is a standard database of RGB-D scene recognition. Similar to previous work Georghiades et al. (2002), we also use a subset of this dataset which contains the 19 major scene categories and 9504 samples in total. Following the author's implementation, We employ pre-trained ResNet-18 as the backbone network for MMTM. The input images are fed into depth and visual block first. Then the rgb and depth features are fused by MMTM before the final prediction. We add CML regularization to the softmax output before and after MMTM fusion process. In our experiment, the squeeze ratio of MMTM Module is set to 16. The dimensionalities of rgb and depth feature are both 512.

Table 8: VRR (%) of test samples (a lower value indicates a better confidence estimation). " λ " indicates the model is not equipped with the proposed regularization ($\lambda = 0$).

Method	CML	NYUD-2	SUN-RGBD
Type III	X ✓ Improve	$\begin{array}{c} 58.09 \pm 4.46 \\ 46.99 \pm 2.89 \\ \bigtriangleup 11.10 \end{array}$	$\begin{array}{c} 57.09 \pm 1.50 \\ 52.56 \pm 3.49 \\ \bigtriangleup 4.53 \end{array}$

Table 9: Accuracy performance comparison of MMTM when some of the modalities is corrupted with color jitter (i.e., randomly change the brightness, contrast, saturation and hue of an image with jitter factor ϵ .).

Dataset	Noise on	CML	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon = 0.5$
NYUD-2	{1}	X ✓ Improve	$\begin{array}{c} 65.72 \pm 0.70 \\ 66.64 \pm 1.22 \\ \triangle 0.92 \end{array}$	$\begin{array}{c} 64.13 \pm 1.78 \\ 65.41 \pm 0.65 \\ \bigtriangleup 1.28 \end{array}$	$\begin{array}{c} 63.79 \pm 1.79 \\ 64.31 \pm 0.92 \\ \bigtriangleup 0.52 \end{array}$	$\begin{array}{c} 60.89 \pm 1.21 \\ 62.26 \pm 1.77 \\ \bigtriangleup 1.37 \end{array}$
	{2}	X ✓ Improve	$\begin{array}{c} 61.34 \pm 0.98 \\ 62.63 \pm 0.60 \\ \bigtriangleup 1.29 \end{array}$	57.98 ± 0.81 57.89 ± 1.56 $\bigtriangledown 0.09$	$\begin{array}{c} 53.98 \pm 2.28 \\ 54.80 \pm 2.90 \\ \bigtriangleup 0.82 \end{array}$	$52.26 \pm 3.23 \\ 52.57 \pm 3.38 \\ \triangle 0.31$
	{1,2}	X ✓ Improve	$\begin{array}{c} 60.43 \pm 0.82 \\ 61.87 \pm 0.93 \\ \bigtriangleup 1.44 \end{array}$	$\begin{array}{c} 55.17 \pm 0.85 \\ 56.24 \pm 2.22 \\ \bigtriangleup 1.07 \end{array}$	$\begin{array}{c} 51.01 \pm 2.64 \\ 51.53 \pm 1.91 \\ \bigtriangleup 0.52 \end{array}$	$\begin{array}{c} 41.52 \pm 4.01 \\ 41.99 \pm 3.37 \\ \bigtriangleup 0.47 \end{array}$
SUN-RGBD	{1}	X ✓ Improve	$\begin{array}{c} 60.72 \pm 0.58 \\ 61.50 \pm 0.59 \\ \bigtriangleup 0.78 \end{array}$	$\begin{array}{c} 58.98 \pm 0.72 \\ 59.95 \pm 0.17 \\ \bigtriangleup 0.97 \end{array}$	$\begin{array}{c} 57.40 \pm 0.75 \\ 57.97 \pm 0.30 \\ \bigtriangleup 0.57 \end{array}$	$\begin{array}{c} 55.68 \pm 0.95 \\ 57.21 \pm 0.32 \\ \bigtriangleup 1.53 \end{array}$
	{2}	X ✓ Improve	$\begin{array}{c} 60.11 \pm 0.24 \\ 59.90 \pm 0.49 \\ \bigtriangledown 0.21 \end{array}$	$58.57 \pm 0.60 \\ 58.44 \pm 0.75 \\ \bigtriangledown 0.13$	$\begin{array}{c} 57.46 \pm 0.69 \\ 57.25 \pm 0.56 \\ \bigtriangledown 0.21 \end{array}$	55.25 ± 1.05 55.34 ± 0.87 -
	{1,2}	X ✓ Improve	$\begin{array}{c} 58.67 \pm 0.42 \\ 58.95 \pm 0.20 \\ \bigtriangleup 0.28 \end{array}$	54.77 ± 0.44 54.73 ± 0.71 -	$\begin{array}{c} 51.66 \pm 0.64 \\ 51.36 \pm 0.66 \\ \bigtriangledown 0.30 \end{array}$	$\begin{array}{c} 45.68 \pm 1.35 \\ 45.99 \pm 1.24 \\ \bigtriangleup 0.31 \end{array}$

G RELATED WORK DETAILS

Uncertainty estimation provides a way for trustworthy prediction (Abdar et al., 2021). Uncertainty can be used as an indicator of whether the predictions given by models are prone to be wrong. Many uncertainty-based models have been proposed in the past decades, such as Bayesian neural networks (Neal, 2012; MacKay, 1992; Denker & LeCun, 1990; Kendall & Gal, 2017), Dropout (Molchanov et al., 2017), and Deep ensembles (Lakshminarayanan et al., 2017; Havasi et al., 2020). Built upon RBF networks, DUQ (van Amersfoort et al., 2020) is able to identify the out-of-distribution samples, which uses distance to represent the prediction uncertainty. Prediction confidence is always referred to in classification models, which expects the predicted class probability to be consistent with the empirical accuracy. Models are frequently overconfident because softmax probabilities are computed with the fast-growing exponential function (Hendrycks & Gimpel, 2017), so many methods focus on smoothing the prediction probabilities distribution, such as Label smoothing (Müller et al., 2019). The recent approach employs the focal loss to calibrate the deep neural networks (Mukhoti et al., 2020). A recent work (Corbière et al., 2019) introduces True Class Probability (TCP) to ensure the low confidence for the failure predictions. Temperature scaling (TS) (Guo et al., 2017) is a well-known post-hoc confidence calibration method, which aims to re-scale the output probability by manipulating the softmax inputs, i.e., the logits.

Incomplete multimodal learning Recently, there have been a wide range of research interests in handling missing modalities for multimodal learning, including imputation-independent (Type I) methods (Zhang et al., 2019) and imputation-dependent (Type II) methods (Mattei & Frellsen, 2019; Wu & Goodman, 2018). Imputation-independent methods have no need to reconstruct the miss-

Dataset	Noise on	CML	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.3$	$\epsilon=0.5$
NYUD-2	{1}	X ✓ Improve	$\begin{array}{c} 64.77 \pm 1.76 \\ 65.26 \pm 1.92 \\ \bigtriangleup 1.49 \end{array}$	$\begin{array}{c} 63.03 \pm 1.92 \\ 63.98 \pm 1.60 \\ \bigtriangleup 0.95 \end{array}$	$\begin{array}{c} 61.50 \pm 2.83 \\ 62.94 \pm 1.97 \\ \bigtriangleup 1.44 \end{array}$	$58.81 \pm 4.05 \\ 59.88 \pm 3.03 \\ \triangle 1.07$
	{2}	X ✓ Improve	$\begin{array}{c} 65.41 \pm 1.27 \\ 66.12 \pm 1.10 \\ \bigtriangleup 1.29 \end{array}$	$\begin{array}{c} 62.17 \pm 1.76 \\ 62.75 \pm 1.26 \\ \bigtriangleup 0.58 \end{array}$	$\begin{array}{c} 59.08 \pm 1.54 \\ 59.79 \pm 2.23 \\ \bigtriangleup 0.71 \end{array}$	$\begin{array}{c} 55.75 \pm 2.75 \\ 55.90 \pm 3.38 \\ \bigtriangleup 0.15 \end{array}$
	{1,2}	X ✓ Improve	$\begin{array}{c} 61.87 \pm 0.82 \\ 63.12 \pm 1.49 \\ \bigtriangleup 1.25 \end{array}$	$\begin{array}{c} 55.60 \pm 2.61 \\ 57.31 \pm 1.58 \\ \bigtriangleup 1.71 \end{array}$	$\begin{array}{c} 48.62 \pm 4.32 \\ 49.51 \pm 2.75 \\ \bigtriangleup 0.89 \end{array}$	$\begin{array}{c} 37.68 \pm 4.94 \\ 37.98 \pm 5.21 \\ \bigtriangleup 0.30 \end{array}$
NYUD-2	{1}	X ✓ Improve	$60.69 \pm 0.65 \\ 61.00 \pm 0.32 \\ riangle 0.31$	$\begin{array}{c} 58.78 \pm 0.95 \\ 59.31 \pm 0.83 \\ \bigtriangleup 0.53 \end{array}$	$\begin{array}{c} 56.84 \pm 1.13 \\ 57.47 \pm 0.62 \\ \bigtriangleup 0.63 \end{array}$	$\begin{array}{c} 53.14 \pm 1.32 \\ 54.77 \pm 1.00 \\ \bigtriangleup 1.63 \end{array}$
	{2}	X ✓ Improve	$\begin{array}{c} 60.93 \pm 0.58 \\ 61.25 \pm 0.59 \\ riangle 0.32 \end{array}$	$59.25 \pm 0.71 \\ 59.19 \pm 0.68 \\ -$	57.55 ± 1.08 57.50 ± 1.27 -	54.81 ± 1.66 54.34 ± 1.93 $\bigtriangledown 0.47$
	{1,2}	X ✓ Improve	$\begin{array}{c} 59.16 \pm 0.88 \\ 59.59 \pm 1.09 \\ \triangle \ 0.43 \end{array}$	$\begin{array}{c} 53.56 \pm 1.51 \\ 54.14 \pm 0.58 \\ \bigtriangleup 0.58 \end{array}$	$\begin{array}{c} 47.22 \pm 2.12 \\ 47.38 \pm 1.47 \\ \triangle \ 0.16 \end{array}$	$\begin{array}{c} 35.90 \pm 2.38 \\ 36.30 \pm 2.39 \\ \triangle \ 0.40 \end{array}$

Table 10: Accuracy performance comparison of MMTM when some of the modalities is corrupted with gaussian noise (i.e., zero mean with varying variance ϵ).

ing modalities and make classification via an uniform representation. For imputation-dependent methods (based on reconstruction), the strategy model can be split into two stages, reconstructing the missing modalities and making classification according to the reconstructed modalities. CPM-Nets (Zhang et al., 2019) is an advanced method (i.e., Type I) which can guarantee the performance by fully exploiting all samples and all modalities to produce structured representation for interpretability, and the method has been extended and deployed into medical domain (Lee & van der Schaar, 2021). MIWAE (Mattei & Frellsen, 2019) is a typical reconstruction model (i.e., Type II) in multimodal learning, whose objective is a lower bound of the likelihood of the observed data that can be tight in the limit of very large computational power.