

# THEORY-INSPIRED TASK-RELEVANT REPRESENTATION LEARNING FOR INCOMPLETE MULTI-VIEW MULTI-LABEL LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Multi-view multi-label learning is commonly hindered by dual data incompleteness, arising from constraints in feature collection and prohibitive annotation costs. To address the intricate yet highly practical challenges and enhance the reliability of representation extraction, heterogeneous feature fusion, and label semantic learning, we propose a Theory-Inspired Task-Relevant Representation Learning method named TITRL. From an information-theoretic standpoint, we identify the sources of view-specific information that interfere with shared representations. By introducing dual-layer constraints on feature exclusivity and label integration, TITRL constructs a general framework for task-relevant information extraction. Besides, through variational derivation, we demonstrate the existence of tractable bounds for the mutual information model that guides the optimization direction. Regarding label semantic learning, we establish flexible relationships between label prototypes by promoting the expression of sample-level label correlations. During the multi-view integration process, TITRL simultaneously incorporates early fusion through distribution information aggregation and late fusion weighted by prediction confidence, which improves the semantic stability while enabling dynamic view quality assessment. Finally, extensive experimental results validate the effectiveness of TITRL against state-of-the-art methods.

## 1 INTRODUCTION

Integrating information from diverse sources enables richer insights and a more holistic understanding of complex problems. By capitalizing on the rapid proliferation of multimodal data, multi-view learning shows strong potential to deliver superior performance across numerous application domains Wen et al. (2022); Fang et al. (2021; 2022). At the same time, multi-label classification, where one instance may correspond to multiple relevant labels, has become increasingly important given the growing annotation demand in the information era. For instance, in remote sensing, an image may be simultaneously tagged with labels such as “forest”, “urban area”, and “water body” Sumbul et al. (2019). Multi-view features provide comprehensive representations of objects, and multiple labels capture their diverse attributes. These characteristics address the limitations of single-view and single-label paradigms in traditional machine learning, aligning with the demands of real-world applications Qin et al. (2025); Tian et al. (2024). By integrating multi-view learning with multi-label classification, a thorough instance depiction and the enriched information for the recognition of multiple labels are attained. Therefore, multi-view multi-label classification (MvMLC) has emerged as a highly promising avenue of research Liu et al. (2025); Zhang et al. (2018).

Existing MvMLC techniques seek to leverage heterogeneous features and predict multiple labels within a unified framework. Representative methods include lrMMC Liu et al. (2015) applying low-rank matrix factorization, and  $EF^2FS$  Hao et al. (2025) based on feature selection. However, many approaches still rely on the assumption that both complete views and full label sets are accessible, which rarely holds in practice. In reality, multi-view data often suffer from missing modalities owing to feature acquisition and processing difficulties. For example, in remote sensing, multispectral imagery may be collected while hyperspectral or LiDAR data are absent due to sensor limitations or high storage requirements Guan et al. (2025). Similarly, multiple annotations are frequently incomplete since labeling cost is expensive, privacy restrictions prevent data sharing, or some categories remain semantically ambiguous. In medical imaging, chest X-rays may contain multiple pathologies but only a subset is annotated, as labeling requires domain expertise and clear diagnostic boundaries are sometimes lacking Sun et al. (2024). The presence of numerous features and labels, coupled

with concurrent data missingness constitutes a widespread challenge, making incomplete multi-view multi-label classification (iMvMLC) particularly complex and urgent to address.

With the advancement of deep learning, various methods based on different network architectures have been applied to address the iMvMLC problem. Nevertheless, these methods still present opportunities for refinement, especially with regard to feature representation extraction, view fusion, and the construction of label semantics. (i) Enhancing information sharing across multiple views is a pivotal factor in both unsupervised clustering tasks Zhou et al. (2024) and supervised classification tasks Chen et al. (2024). DICNet Liu et al. (2023b) and LMVCAT Liu et al. (2023c) capture shared representations by utilizing cross-view interaction mechanisms. However, these methods fail to account for the disruptions caused by view-specific information. As a result, redundancy and noise are inadvertently incorporated into the shared representations, diminishing their purity and increasing the risk of misguiding the classification process. Although SIP Liu et al. (2024b) is a method for minimizing non-shared information and maintaining feature validity, it does not integrate label information to guide representation extraction, which leads to uncertainty about the practicality of the obtained common information. (ii) Previous approaches, such as DIMC Wen et al. (2023) and AIMNet Liu et al. (2024a), Wen et al. (2023); Liu et al. (2024a) have largely focused on feature-level weighting for view fusion. Nevertheless, without leveraging classification confidence as a weighting signal, such methods often fail to capture discriminative information from each view. As the number of categories increases, it becomes crucial to identify pertinent information for predicting each category, which underscores the need for label-specific feature selection. (iii) Learning multi-label semantics necessitates modeling label relationships. Methods like MTD Liu et al. (2023a), which treat multi-label learning as separate binary classifications, are inherently limited in achieving optimal performance. Moreover, label correlations cannot be regarded as fixed pairwise measures, as assumed in traditional methods. The realization of label correlations often fluctuates between different instances Si et al. (2023). For example, in a movie recommendation system Li et al. (2025), the correlation between "action" and "adventure" genres may be stronger for some users, while weaker for others, depending on individual preferences.

To address these problems, we propose a Theory-Inspired Task-Relevant Representation Learning framework named TITRL. The motivation behind TITRL is to enhance the purity of shared representations, improve the effectiveness of view fusion, and delicately capture the multi-label correlation semantics. We begin by leveraging mutual information-based semantic interaction and theoretically establishing a dual-layer constraint framework at the levels of feature and category. Guided by the principle of mitigating view-specific noise that adversely affects representation extraction and downstream prediction, we disentangle the view-specific mixtures that indicate the negative influence of each view on label recognition. Besides, we obtain tractable bounds for the mutual information model through variational derivation, which serves as the training loss to guide the extraction of common information. Regarding view fusion, we initially employ a distribution-aware blending strategy to derive the distribution parameters of the integrated shared information, which not only aids in selecting views with stable statistical properties but also facilitates coherent posterior distribution inference. After constructing the prototype representation for each label, the pseudo-labels are generated by leveraging the interactions between view representations and these prototypes. Subsequently, we perform a confidence-based late fusion by utilizing the pseudo-labels derived from the remaining views after removing each individual view, along with the prediction from all views. The process aims to mitigate the view-specific interference while retaining the most informative insights that contribute to label prediction. Finally, to accurately model label correlations, we focus on maximizing the similarity between the positive label prototypes of each sample and its shared representation. This approach promotes the learning of a sample-specific correlation structure, which enables flexible utilization of label dependencies to improve classification performance. The main contributions of our TITRL are summarized as follows:

- We propose a general framework for multi-view shared representation extraction, applying constraints at both the feature and label levels. Moreover, we theoretically establish the optimization direction of the model and derive the variational bound to guide the training process.
- TITRL simultaneously considers the statistical properties of representation extraction and the confidence of label prediction in view fusion. Additionally, TITRL proposes a flexible approach to represent label correlations, which focuses on the diverse manifestation patterns across samples.
- Extensive experimental results across a range of public datasets and varying degrees of data missingness demonstrate the effectiveness and robustness of our method.

## 2 METHOD

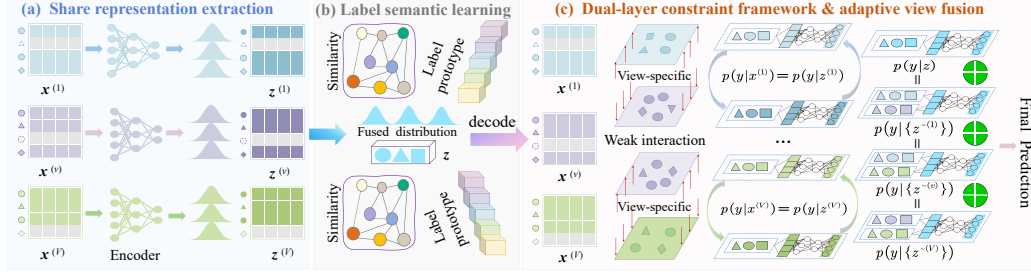


Figure 1: The main framework of our proposed TITRL. Different shapes signify different samples.

### 2.1 PROBLEM DEFINITION

Consider a dataset consisting of  $n$  labeled instances, represented as  $(\{\mathbf{x}^{(v)}\}_{v=1}^V, \mathbf{y})$ , where each sample is observed from  $V$  distinct views. Specifically, the  $v$ -th view of any sample is denoted as  $\mathbf{x}^{(v)} \in \mathbb{R}^{d_v}$ , while the associated label  $\mathbf{y} \in \{0, 1\}^c$  corresponds to  $c$  categories. Additionally, we define  $\mathcal{V}$  ( $|\mathcal{V}| \leq V$ ) as the set of observed views. Thus, the available multi-view data can be expressed as  $\{\mathbf{x}^{(v)}\}_{v \in \mathcal{V}}$  (abbreviated as  $\{\mathbf{x}\}$ ). Moreover, let  $\mathcal{U}$  represent the set of known tags, where  $|\mathcal{U}| \leq c$ . The goal is to design an end-to-end neural network capable of performing classification tasks on incomplete multi-view partial multi-label data.

### 2.2 TASK-RELEVANT REPRESENTATION LEARNING UNDER A DUAL-LAYER CONSTRAINT FRAMEWORK

Enhancing cross-view information interaction in multi-view learning has consistently been a crucial driver of improved classification performance. Moreover, prior researches Federici et al. (2020) have demonstrated that integrating the common information across all views and reducing the redundant information introduced by view-specific factors is sufficient to accomplish all prediction tasks. For example, in facial recognition Liu et al. (2023d), images from different views capture shared facial features, with the frontal view revealing details of the eyes and nose, and the side view presenting the contours. However, some views may introduce disruptive factors, including lighting variations, background clutter, or excessive emphasis on minor details, which can disrupt model performance. By integrating shared features and removing noisy information, the model can achieve more accurate recognition. Given an initial shared representation  $\mathbf{z}^{(v)} \in \mathbb{R}^d$  for each view, the unified representation  $\mathbf{z} \in \mathbb{R}^d$  is obtained by fusing them. To guarantee that the shared representation captures the common information across all views, it is crucial for the semantics of  $\mathbf{z}$  to encompass the relevant information from original views as much as possible. This objective introduces the requirement of optimizing the mutual information interactions between  $\mathbf{z}$  and each individual view to their fullest extent, i.e.,  $\max \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} I(\mathbf{x}^{(v)}; \mathbf{z})$ . Moreover, the detrimental redundancy arising from the distinct information inherent to each view should be meticulously minimized, which necessitates that the derived representation primarily conveys the shared components, while effectively attenuating the noise caused by view-specific characteristics to the absolute minimum. Thus, the representation  $\mathbf{z}^{(v)}$  ought to be distanced from the view-specific information from other perspectives, with the aim of minimizing  $I(\{\mathbf{x}^{(v)}\}; \mathbf{z}^{(v)} | \mathbf{x}^{(v)})$ , where  $\{\{\mathbf{x}^{(v)}\}, \mathbf{x}^{(v)}\} = \{\mathbf{x}\}$ . Next, due to the scalability of information transfer, we can derive the following upper bound to guide information separation:

$$\sum_{v \in \mathcal{V}} I(\{\mathbf{x}^{(v)}\}; \mathbf{z}^{(v)} | \mathbf{x}^{(v)}) \leq \sum_{v \in \mathcal{V}} I(\{\mathbf{x}^{(v)}\}; \mathbf{z} | \mathbf{x}^{(v)}). \quad (1)$$

We have concentrated on the suppression of view-specific redundancy at the feature level. However, it remains uncertain whether these representations are directly applicable to downstream classification as label information is not integrated. Therefore, it is crucial to incorporate task-specific knowledge to steer the unification of these features toward enhancing classification performance. Foremost, it is imperative to prevent information degradation by ensuring that the extracted representations preserve the mutual information between the original features and their corresponding labels. This requirement imposes the exact equivalence between  $I(\mathbf{x}^{(v)}; \mathbf{y})$  and  $I(\mathbf{z}^{(v)}; \mathbf{y})$ :

$$\min \sum_{v \in \mathcal{V}} (I(\mathbf{x}^{(v)}; \mathbf{y}) - I(\mathbf{z}^{(v)}; \mathbf{y})). \quad (2)$$

In addition, another crucial consideration lies in ensuring the extracted information is solely label-relevant and devoid of any admixed noise. In this regard, by isolating the distinctive impact of  $\mathbf{z}^{(v)}$  within the task-relevant components, we obtain the following expression:

$$I(\mathbf{y}; \mathbf{z}^{(v)}) = \underbrace{\sum_{j=1, j \neq v}^V I(\mathbf{y}; \{\mathbf{z}^{(j)}\} | \mathbf{z}^{(j)})}_{\text{shared } I_v^s} + \underbrace{I(\mathbf{y}; \{\mathbf{z}\}) + I(\mathbf{y}; \mathbf{z}^{(v)} | \{\mathbf{z}^{(v)}\})}_{\text{view-specific}}, \quad (3)$$

where the preceding term is referred to as shared information, as each of its components encapsulates associative information contributed collectively by multiple views toward the label. Then, our optimization goal is to achieve cleaner feature extraction by controlling task-irrelevant information, reduce misclassifications caused by view-specific redundancy, and emphasize the collaborative discriminative power of all useful signals from multi-view. Since the shared information term  $I_v^s$  consists of multiple components and cannot be directly optimized, we substitute it with its upper bound  $I(\mathbf{y}; \mathbf{z}^{(v)})$  based on its optimization direction. Therefore, under the dual-layer constraints at both the feature and category levels, the model for acquiring shared representations is obtained:

$$\min \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \left( \underbrace{-I(\mathbf{x}^{(v)}; \mathbf{z}) + I(\{\mathbf{x}^{(v)}\}; \mathbf{z} | \mathbf{x}^{(v)})}_{\text{feature-level}} + \underbrace{I(\mathbf{x}^{(v)}; \mathbf{y}) - I(\mathbf{z}^{(v)}; \mathbf{y}) - I(\mathbf{y}; \mathbf{z}^{(v)}) + I(\mathbf{y}; \mathbf{z}^{(v)} | \{\mathbf{z}^{(v)}\})}_{\text{category-level}} \right). \quad (4)$$

Due to the intractability of computing mutual information in high-dimensional spaces, we derive its bound that allows for reliable estimation to facilitate the optimization of model (4). For the first term  $I(\mathbf{x}^{(v)}; \mathbf{z})$ , its lower bound is typically expressed via a reconstruction loss, where  $\mathbf{x}^{(v)}$  is decoded through the decoder  $q^v(\mathbf{x}^{(v)} | \mathbf{z})$  to ensure the faithful preservation of the original view:

$$I(\mathbf{x}^{(v)}; \mathbf{z}) \geq \mathbb{E}_{p(\mathbf{x}^{(v)}; \mathbf{z})} [\log q^v(\mathbf{x}^{(v)} | \mathbf{z})] = \mathbb{E}_{\{\mathbf{x}\} \sim p(\{\mathbf{x}\})} \left[ \int p(\mathbf{z} | \{\mathbf{x}\}) \log q^v(\mathbf{x}^{(v)} | \mathbf{z}) d\mathbf{z} \right]. \quad (5)$$

Next, based on the definition of mutual information, the expansion for the second term is derived:

$$\begin{aligned} I(\{\mathbf{x}^{(v)}\}; \mathbf{z} | \mathbf{x}^{(v)}) &= \mathbb{E}_{p(\{\mathbf{x}^{(v)}\}, \mathbf{z}, \mathbf{x}^{(v)})} \left[ \log \frac{p(\{\mathbf{x}^{(v)}\}, \mathbf{z} | \mathbf{x}^{(v)})}{p(\{\mathbf{x}^{(v)}\} | \mathbf{x}^{(v)}) p(\mathbf{z} | \mathbf{x}^{(v)})} \right] \\ &= \int \int p(\{\mathbf{x}\}, \mathbf{z}) \log \frac{p(\mathbf{z} | \{\mathbf{x}\})}{p(\mathbf{z} | \mathbf{x}^{(v)})} d\{\mathbf{x}\} d\mathbf{z}. \end{aligned} \quad (6)$$

Since the distribution  $p(\mathbf{z} | \{\mathbf{x}\})$  is difficult to obtain explicitly, we approximate it using a stochastic variational distribution  $g^v(\mathbf{z} | \mathbf{x}^{(v)})$ . Then, we can obtain the following transformation:

$$\begin{aligned} I(\{\mathbf{x}^{(v)}\}; \mathbf{z} | \mathbf{x}^{(v)}) &= \int \int p(\{\mathbf{x}\}, \mathbf{z}) \log \frac{p(\mathbf{z} | \{\mathbf{x}\})}{g^v(\mathbf{z} | \mathbf{x}^{(v)})} d\{\mathbf{x}\} d\mathbf{z} - \int p(\mathbf{x}^{(v)}) D_{KL}(p(\mathbf{z} | \mathbf{x}^{(v)}) || g^v(\mathbf{z} | \mathbf{x}^{(v)})) d\mathbf{x}^{(v)}, \end{aligned} \quad (7)$$

where  $D_{KL}(\cdot || \cdot)$  denotes the non-negative Kullback-Leibler divergence. Thus, the variational upper bound for  $I(\{\mathbf{x}^{(v)}\}; \mathbf{z} | \mathbf{x}^{(v)})$  can be established:

$$I(\{\mathbf{x}^{(v)}\}; \mathbf{z} | \mathbf{x}^{(v)}) \leq \mathbb{E}_{\{\mathbf{x}\} \sim p(\{\mathbf{x}\})} [D_{KL}(p(\mathbf{z} | \{\mathbf{x}\}) || g^v(\mathbf{z} | \mathbf{x}^{(v)}))]. \quad (8)$$

In summary, the trainable loss subject to the feature-level constraint is given by:

$$\mathcal{L}_f = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \left[ -\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z} | \{\mathbf{x}\})} \log q^v(\mathbf{x}^{(v)} | \mathbf{z}) + D_{KL}(p(\mathbf{z} | \{\mathbf{x}\}) || g^v(\mathbf{z} | \mathbf{x}^{(v)})) \right]. \quad (9)$$

Under category-level constraints, the minimization of  $I(\mathbf{x}^{(v)}; \mathbf{y}) - I(\mathbf{z}^{(v)}; \mathbf{y})$  is functionally equivalent to restricting  $H(\mathbf{y} | \mathbf{z}^{(v)}) - H(\mathbf{y} | \mathbf{x}^{(v)})$ , where  $H(\cdot)$  denotes the Shannon entropy. Since the disparity in entropy is characterized by the divergence between distributions, the constraint objective naturally transitions to:

$$\min \sum_{v \in \mathcal{V}} D_{KL} \left( p(\mathbf{y}|\mathbf{z}^{(v)}) || p(\mathbf{y}|\mathbf{x}^{(v)}) \right) \quad (10)$$

Regarding the latter part of Model (4), the view-specific information is decomposed as follows:

$$\begin{aligned} I(\mathbf{y}; \mathbf{z}_i | \{\mathbf{z}^{\sim(v)}\}) &= H(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) - H(\mathbf{y} | \{\mathbf{z}\}) \\ &= - \int p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \log p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) d\mathbf{y} + \int p(\mathbf{y} | \{\mathbf{z}\}) \log p(\mathbf{y} | \{\mathbf{z}\}) d\mathbf{y}, \end{aligned} \quad (11)$$

Through term augmentation and subsequent expansion in logarithmic operations, we have

$$\begin{aligned} I(\mathbf{y}; \mathbf{z}^{(v)} | \{\mathbf{z}^{\sim(v)}\}) &= - \int p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \log \left[ \frac{p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})}{p(\mathbf{y} | \{\mathbf{z}\})} \right] d\mathbf{y} - \int p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \log p(\mathbf{y} | \{\mathbf{z}\}) d\mathbf{y} \\ &\quad + \int p(\mathbf{y} | \{\mathbf{z}\}) \log \left[ \frac{p(\mathbf{y} | \{\mathbf{z}\})}{p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})} \right] d\mathbf{y} + \int p(\mathbf{y} | \{\mathbf{z}\}) \log p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) d\mathbf{y}. \\ &\leq D_{KL} \left( p(\mathbf{y} | \{\mathbf{z}\}) || p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \right) + H \left( p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}), p(\mathbf{y} | \{\mathbf{z}\}) \right). \end{aligned} \quad (12)$$

Owing to the congruent optimization objective of aligning  $p(\mathbf{y} | \{\mathbf{z}\})$  and  $p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})$ , we exclusively adopt  $D_{KL} [p(\mathbf{y} | \{\mathbf{z}\}) || p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})]$  as the minimization target. Meanwhile, the existence of Eq. (3) enables the maximal elimination of view-specific noise to accentuate the label-related consensus information  $I_v^s$ . Then, the objective function guided by the category-level constraint is formulated as

$$\mathcal{L}_c = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \left( D_{KL} \left( p(\mathbf{y} | \mathbf{z}^{(v)}) || p(\mathbf{y} | \mathbf{x}^{(v)}) \right) + D_{KL} \left( p(\mathbf{y} | \{\mathbf{z}\}) || p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \right) \right). \quad (13)$$

### 2.3 ADAPTIVE VIEW FUSION AND LABEL REPRESENTATION LEARNING

The integration of view representations constitutes a critical challenge in multi-view learning. Given that variational inference optimizes the distribution of each view, we leverage the progressively refined distribution information to facilitate view fusion. Within the network architecture, each view is processed through two encoders to estimate the latent distribution of its shared representations. Specifically, the distribution is modeled as  $p(\mathbf{z}^{(v)} | \mathbf{x}^{(v)}) := \mathcal{N}(f_\mu^v(\mathbf{x}^{(v)}), f_{\sigma^2}^v(\mathbf{x}^{(v)})\mathbf{I})$ , where  $f_\mu^v$  and  $f_{\sigma^2}^v$  are the mean and variance encoders. To ensure that the shared feature incorporates information from all views and benefits from the greater stability of representations with lower variance, we adopt the product-of-experts (PoE) framework (Hinton, 2002) with the one-vote property to perform weighted fusion of the distribution parameters across views:

$$\begin{cases} \mu_s = \frac{\sum_{v \in \mathcal{V}} f_\mu^v(\mathbf{x}^{(v)}) \frac{1}{f_{\sigma^2}^v(\mathbf{x}^{(v)})}}{\sum_{v \in \mathcal{V}} \frac{1}{f_{\sigma^2}^v(\mathbf{x}^{(v)})} + 1} \\ \sigma_s^2 = \frac{1}{\sum_{v \in \mathcal{V}} \frac{1}{f_{\sigma^2}^v(\mathbf{x}^{(v)})} + 1}. \end{cases} \quad (14)$$

Then, we employ the reparameterization trick to sample  $S$  times from the distribution:

$$\mathbf{z} = \frac{1}{S} \sum_{i=1}^S (\mu_s + \sigma_s \odot \delta^i), \quad (15)$$

where  $\delta^i \in \mathbb{R}^d$  denotes the  $i$ -th sampling from the standard Gaussian distribution and  $\odot$  indicates element-wise multiplication. The representations  $\{\mathbf{z}^{(v)}\}_{v=1}^V$  extracted from each view are also sampled from their respective distributions following Eq. (15). During view fusion, it is essential to not only account for the aggregation of representation information but also to incorporate the impact

of label information. Since multiple labels are typically encoded as one-hot vectors, which lacks the flexibility to capture label semantics, particularly in scenarios with missing labels. To address this, we adopt a data-driven approach to introduce label prototypes, ensuring that the semantic information carried by these prototypes is closely aligned with the ground truth labels. In order to explicitly model label expressions, we employ stochastic encoders to fit the underlying distribution  $\mathcal{N}(\mu_i, \sigma_i^2 \mathbf{I})$  for each label prototype, where  $\mu_i$  and  $\sigma_i^2$  are the  $d$ -dimensional mean and variance outputs, respectively, produced by the encoders  $h_\mu(\mathbf{b}_i)$  and  $h_{\sigma^2}(\mathbf{b}_i)$ .  $\mathbf{b}_i \in \mathbb{R}^C$  serves as a learnable embedding corresponding to the  $i$ -th class, which is initialized as a one-hot vector with the  $i$ -th entry is 1. After obtaining  $\{\mathbf{l}_i\}_{i=1}^C$  through stochastic sampling, it is necessary to capture the intricate correlations between these label representations, which forms a crucial determinant in enhancing the performance of multi-label classification. Considering that the manifestation of label correlations differs across samples, we adopt a nuanced approach that centers on instance-level relevance to strengthen the similarity between the cross-view representation of each sample and the label attributes it possesses. Specifically, we sample the shared feature  $\mathbf{z}$  according to Eq. (15), with its associated known label prototype being  $\{\mathbf{l}_i | i \in \mathcal{U}\}$ . By using the cosine similarity as the criterion, the alignment loss designed to capture label correlations is as follows:

$$\mathcal{L}_a = \frac{1}{|\mathcal{U}|} \sum_{i \in \mathcal{U}} \frac{\langle \mathbf{z} \cdot \mathbf{l}_i \rangle}{\|\mathbf{z}\| \|\mathbf{l}_i\|}. \quad (16)$$

By optimizing loss  $\mathcal{L}_a$ , we refine the mapping semantic between features and labels, while simultaneously highlighting the associations among label prototypes, which are tailored for application to each individual sample. Subsequently, during label prediction, it is important to synthesize the generalized information from multiple views with the semantic representations of individual categories. When these information exhibit coherence, it becomes feasible to infer that the sample contains the relevant labels. To this end, we utilize a neural network to adaptively gauge the degree of similarity between view representations and category embeddings:

$$\mathbf{p}_i^0 = \omega(g_c(\mathbf{z} \oplus \mathbf{l}_i)), \quad (17)$$

where  $g_c$  is a fully connected layer,  $\oplus$  denotes concatenation operation and  $\sigma_S$  is the Sigmoid activation function. The derivation of  $\mathbf{p}_i^0$  solely relies on the shared representation  $\mathbf{z}$  resulting from the fusion of feature information. To further refine the integration of effective multimodal information, we incorporate multi-label semantic information into the fusion process. To achieve this, we propose a label-guided post-view fusion framework, where  $\mathbf{p}^0$  and the label distributions  $p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})$  obtained from the exclusion of each view are adaptively merged. This strategy is designed to mitigate the adverse effects of heterogeneous views on label recognition while preserving the most discriminative feature information, thereby improving the reliability of the prediction outcome. Then, we utilize the computed result  $\mathcal{L}_{con} = \frac{1}{|\mathcal{U}|} \sum_{i \in \mathcal{U}} (\mathbf{p}_i^2 + (1 - \mathbf{p}_i)^2)$  of the predicted label distribution as its confidence measure. Besides, we derive  $\mathcal{L}_{con}^0$  based on  $\mathbf{p}^0$ , and calculate  $\mathcal{L}_{con}^{(v)}$  from  $p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})$ . The formulation of  $\mathcal{L}_{con}$  indicates that the value of 0.5 serves as the classification boundary. Scores significantly exceeding 0.5 indicate a stronger tendency toward positive labels, while those substantially below 0.5 reflect an increased likelihood of negative assignment. Therefore, by employing  $\mathcal{L}_{con}$  as the weighting factor for late fusion, we can obtain the enhanced result as the final prediction:

$$\mathbf{p}_i^t = \sum_{v \in \mathcal{V}} \mathcal{L}_{con}^{(v)} p(\mathbf{y}_i | \{\mathbf{z}^{\sim(v)}\}) + \mathcal{L}_{con}^{(0)} \mathbf{p}_i^0, \quad (18)$$

where all weighting coefficients  $\mathcal{L}_{con}^{(v)} (0 \leq v \leq V)$  are normalized in advance. To enhance the classification discriminability and reinforce the interaction term  $I(\mathbf{y}; \mathbf{z}^{(v)})$  in model (4), we employ the following cross-entropy loss:

$$\mathcal{L}_{BCE} = \frac{1}{|\mathcal{U}|} \sum_{i \in \mathcal{U}} [\mathbf{y}_i \log \mathbf{p}_i + (1 - \mathbf{y}_i) \log (1 - \mathbf{p}_i)]. \quad (19)$$

The classification loss in our method is the aggregation of four distinct components, with one arising from the final prediction and the remaining three emanating from pseudo-predictions  $p(\mathbf{y}_i | \{\mathbf{z}^{\sim(v)}\})$ ,  $p(\mathbf{y}_i | \mathbf{x}^{(v)})$ , and  $p(\mathbf{y}_i | \mathbf{z}^{(v)})$ , which collectively constitutes the overall loss  $\mathcal{L}_{BCE}^t$ . Thus, the total training loss of TITRR is as below:

$$\mathcal{L} = \mathcal{L}_{BCE}^t + \mathcal{L}_a + \lambda_1 \mathcal{L}_f + \lambda_2 \mathcal{L}_c, \quad (20)$$

where  $\lambda_1$  and  $\lambda_2$  govern the trade-off between the empirical values and impacts of different losses.

### 3 EXPERIMENTS

#### 3.1 DATASETS AND METRICS

In our experiments, we employ six widely used multi-view multi-label datasets to evaluate the effectiveness of our method, i.e., Corel 5k Duygulu et al. (2002), ESPGame Ahn & Dabbish (2004), IAPRTC12 Grubinger et al. (2006), Mirflickr Huiskes & Lew (2008), Pascal07 Everingham et al. (2010), and OBJECT Hao et al. (2024). Following the evaluation protocols in Liu et al. (2023b); Wen et al. (2023), we adopt the following six metrics to form a comprehensive assessment framework, i.e., Hamming Loss (HL), Ranking Loss (RL), OneError (OE), Coverage (Cov), Average Precision (AP), and Area Under Curve (AUC). For clarity in comparison, we report 1-HL, 1-OE, 1-Cov, and 1-RL, where higher values consistently indicate better performance.

#### 3.2 COMPARISON METHODS

To assess the performance of our method, we compare it with nine state-of-the-art approaches, i.e., AIMNet Liu et al. (2024a), DICNet Liu et al. (2023b), DIMC Wen et al. (2023), iMVWL Tan et al. (2018), LMVCAT Liu et al. (2023c), MTD Liu et al. (2023a), SIP Liu et al. (2024b), LVSL Zhao et al. (2022), and DM2L Ma & Chen (2021). The first seven methods are capable of simultaneously handling missing views and labels. Since LVSL cannot directly process incomplete data, we impute missing views using the mean of available instances and fill absent labels with zeros. DM2L is a kernel-based nonlinear method for incomplete multi-label learning. Thus, we concatenate all recovered views into a single representation to apply DM2L. All hyperparameters of the compared methods are set according to the recommended configurations in their original implementations, ensuring a fair and reproducible comparison.

#### 3.3 IMPLEMENTATION DETAILS

To simulate partial view scenarios, a proportion of instances determined by the Partial Example Ratio (PER) are randomly masked in each view, while ensuring each sample retains at least one complete view. For weak supervision, label omissions are applied to both positive and negative tags according to the Label Missing Ratio (LMR). Incomplete data construction is repeated multiple times to mitigate randomness. Datasets are split into training, validation, and test sets with a 7:1:2 ratio. Our method is implemented in PyTorch and trained on an NVIDIA GeForce RTX 4090 GPU.

#### 3.4 EXPERIMENTAL RESULTS AND ANALYSIS

To rigorously evaluate the effectiveness of TITRL in handling absent views and incomplete labels, we conduct extensive comparative experiments against nine representative algorithms across six benchmark datasets under varying levels of data sparsity. Specifically, the proportions of missing views (PER) and labels (LMR) are set to  $\{30\%, 50\%, 70\%, 90\%\}$ . The results in terms of the mean and standard deviation at PER=50% and LMR=50% are summarized in Table 1, along with the average ranking across six evaluation metrics to provide an aggregated performance assessment. In addition, Fig. 2 visualizes how AP evolves as the missing proportion increases, while Fig. 3 presents radar plots that jointly capture multi-metric performance distribution at PER=90% and LMR=90%. These results collectively provide a holistic evaluation of predictive accuracy and model robustness.

From the comparison results, several important observations can be drawn: (i) TITRL consistently secures the best results across almost all datasets and metrics. For instance, on Corel5k, TITRL achieves an AP score of 0.432, outperforming SIP (0.414) and MTD (0.410), with similar margins observed on other datasets. Besides, TITRL maintains its superiority on large-scale datasets under data missingness like ESPGame and Mirflickr, which underscores its scalability and resilience. (ii) Drawn from Fig. 2, we can find that while competing methods suffer steep performance degradation when missing ratios achieve a high level, TITRL continues to exhibit considerable performance. For example, when PER=70%, the performance of all comparison methods stays below 0.36, while TITRL surpasses 0.4 by a certain margin. Although our method still outperforms others under conditions of severe label missingness, the performance gap is prominently reflected in the presence of feature absent. This further highlights the critical importance of multi-view representation learning, a role that our method is well equipped to fulfill. As depicted in the radar chart of Fig. 3, it is evident that our method consistently occupies the outermost boundary, which indicates that TITRL stands out even under highly challenging conditions across various evaluation perspectives. Consequently, our method demonstrates strong robustness in addressing the problem of data incompleteness and

shows significant potential for broader adoption. (iii) As evidenced by Table 1, our method consistently maintains the top position, while the rankings of alternative approaches remain volatile, which shows the exceptional performance stability of our method. Against top competing methods SIP and MTD, our approach also demonstrates superiority, which underscores the pivotal role of introducing label integration strategies in the feature extraction process and fine-grained characterization of label semantics. Compared with the deep learning-based methods AIMNet, DICNet and DIMC, the substantial advantage of TITRL further reveals the importance of jointly considering the view property and label information during the fusion process.

Table 1: Experimental results of nine methods on the six datasets with 50% PER and 50% LMR. ‘AVE’ refers to the mean ranking of the corresponding method across all six metrics. The best and second best results are highlighted in red and blue, respectively.

DATA	METRIC	AIMNet	DICNet	DIMC	DM2L	iMVWL	LMVCAT	LVSL	MTD	SIP	TITRL
COR	1-HL	0.988 <sub>0.000</sub>	0.987 <sub>0.000</sub>	0.987 <sub>0.000</sub>	0.987 <sub>0.000</sub>	0.978 <sub>0.000</sub>	0.986 <sub>0.000</sub>	0.987 <sub>0.000</sub>	0.988 <sub>0.000</sub>	0.988 <sub>0.000</sub>	0.988 <sub>0.000</sub>
	1-OE	0.478 <sub>0.011</sub>	0.460 <sub>0.012</sub>	0.446 <sub>0.009</sub>	0.378 <sub>0.014</sub>	0.308 <sub>0.017</sub>	0.448 <sub>0.011</sub>	0.353 <sub>0.017</sub>	0.492 <sub>0.011</sub>	0.492 <sub>0.014</sub>	0.509 <sub>0.014</sub>
	1-Cov	0.766 <sub>0.004</sub>	0.726 <sub>0.007</sub>	0.709 <sub>0.008</sub>	0.640 <sub>0.007</sub>	0.701 <sub>0.003</sub>	0.720 <sub>0.006</sub>	0.720 <sub>0.005</sub>	0.754 <sub>0.005</sub>	0.781 <sub>0.004</sub>	0.795 <sub>0.006</sub>
	1-RL	0.900 <sub>0.002</sub>	0.881 <sub>0.004</sub>	0.874 <sub>0.004</sub>	0.843 <sub>0.004</sub>	0.864 <sub>0.002</sub>	0.876 <sub>0.004</sub>	0.879 <sub>0.002</sub>	0.893 <sub>0.004</sub>	0.908 <sub>0.003</sub>	0.914 <sub>0.003</sub>
	AP	0.404 <sub>0.005</sub>	0.381 <sub>0.006</sub>	0.370 <sub>0.005</sub>	0.318 <sub>0.005</sub>	0.281 <sub>0.005</sub>	0.379 <sub>0.006</sub>	0.311 <sub>0.005</sub>	0.410 <sub>0.007</sub>	0.414 <sub>0.006</sub>	0.432 <sub>0.007</sub>
	AUC	0.903 <sub>0.002</sub>	0.883 <sub>0.004</sub>	0.877 <sub>0.004</sub>	0.846 <sub>0.004</sub>	0.867 <sub>0.002</sub>	0.879 <sub>0.003</sub>	0.882 <sub>0.002</sub>	0.896 <sub>0.003</sub>	0.910 <sub>0.002</sub>	0.916 <sub>0.002</sub>
	AVE	3.5	5.0	7.2	9.0	9.5	6.8	7.3	3.2	2.2	1.0
ESP	1-HL	0.983 <sub>0.000</sub>	0.983 <sub>0.000</sub>	0.983 <sub>0.000</sub>	0.983 <sub>0.000</sub>	0.972 <sub>0.000</sub>	0.982 <sub>0.000</sub>	0.983 <sub>0.000</sub>	0.983 <sub>0.000</sub>	0.983 <sub>0.000</sub>	0.983 <sub>0.000</sub>
	1-OE	0.442 <sub>0.006</sub>	0.440 <sub>0.009</sub>	0.431 <sub>0.009</sub>	0.302 <sub>0.008</sub>	0.343 <sub>0.010</sub>	0.432 <sub>0.006</sub>	0.365 <sub>0.006</sub>	0.452 <sub>0.007</sub>	0.450 <sub>0.006</sub>	0.481 <sub>0.006</sub>
	1-Cov	0.621 <sub>0.003</sub>	0.601 <sub>0.003</sub>	0.586 <sub>0.004</sub>	0.532 <sub>0.003</sub>	0.548 <sub>0.004</sub>	0.587 <sub>0.003</sub>	0.578 <sub>0.002</sub>	0.617 <sub>0.004</sub>	0.622 <sub>0.004</sub>	0.631 <sub>0.008</sub>
	1-RL	0.845 <sub>0.002</sub>	0.836 <sub>0.002</sub>	0.830 <sub>0.002</sub>	0.804 <sub>0.002</sub>	0.807 <sub>0.002</sub>	0.827 <sub>0.002</sub>	0.829 <sub>0.001</sub>	0.843 <sub>0.002</sub>	0.847 <sub>0.002</sub>	0.852 <sub>0.004</sub>
	AP	0.306 <sub>0.003</sub>	0.300 <sub>0.003</sub>	0.294 <sub>0.003</sub>	0.229 <sub>0.003</sub>	0.243 <sub>0.004</sub>	0.293 <sub>0.003</sub>	0.266 <sub>0.003</sub>	0.309 <sub>0.003</sub>	0.309 <sub>0.004</sub>	0.339 <sub>0.003</sub>
	AUC	0.850 <sub>0.001</sub>	0.841 <sub>0.002</sub>	0.835 <sub>0.002</sub>	0.808 <sub>0.001</sub>	0.813 <sub>0.002</sub>	0.832 <sub>0.001</sub>	0.834 <sub>0.001</sub>	0.847 <sub>0.002</sub>	0.851 <sub>0.002</sub>	0.855 <sub>0.003</sub>
	AVE	3.7	4.5	5.7	9.7	9.2	7.3	7.2	3.5	2.3	1.0
IAP	1-HL	0.981 <sub>0.000</sub>	0.981 <sub>0.000</sub>	0.981 <sub>0.000</sub>	0.980 <sub>0.000</sub>	0.969 <sub>0.000</sub>	0.980 <sub>0.000</sub>	0.981 <sub>0.000</sub>	0.981 <sub>0.000</sub>	0.981 <sub>0.000</sub>	0.982 <sub>0.000</sub>
	1-OE	0.457 <sub>0.008</sub>	0.464 <sub>0.008</sub>	0.454 <sub>0.006</sub>	0.378 <sub>0.008</sub>	0.351 <sub>0.008</sub>	0.433 <sub>0.009</sub>	0.377 <sub>0.007</sub>	0.479 <sub>0.007</sub>	0.459 <sub>0.005</sub>	0.508 <sub>0.008</sub>
	1-Cov	0.675 <sub>0.004</sub>	0.649 <sub>0.005</sub>	0.630 <sub>0.005</sub>	0.555 <sub>0.005</sub>	0.565 <sub>0.004</sub>	0.646 <sub>0.004</sub>	0.605 <sub>0.004</sub>	0.670 <sub>0.004</sub>	0.678 <sub>0.003</sub>	0.693 <sub>0.006</sub>
	1-RL	0.884 <sub>0.001</sub>	0.874 <sub>0.002</sub>	0.868 <sub>0.002</sub>	0.837 <sub>0.002</sub>	0.833 <sub>0.002</sub>	0.868 <sub>0.002</sub>	0.857 <sub>0.002</sub>	0.882 <sub>0.002</sub>	0.886 <sub>0.001</sub>	0.893 <sub>0.002</sub>
	AP	0.329 <sub>0.003</sub>	0.326 <sub>0.003</sub>	0.318 <sub>0.002</sub>	0.254 <sub>0.002</sub>	0.236 <sub>0.002</sub>	0.313 <sub>0.004</sub>	0.262 <sub>0.002</sub>	0.340 <sub>0.002</sub>	0.331 <sub>0.002</sub>	0.377 <sub>0.004</sub>
	AUC	0.885 <sub>0.001</sub>	0.876 <sub>0.002</sub>	0.870 <sub>0.001</sub>	0.838 <sub>0.001</sub>	0.835 <sub>0.001</sub>	0.870 <sub>0.002</sub>	0.859 <sub>0.001</sub>	0.883 <sub>0.002</sub>	0.887 <sub>0.001</sub>	0.894 <sub>0.002</sub>
	AVE	4.0	4.3	6.0	8.8	9.8	6.8	8.0	3.0	2.8	1.0
MIR	1-HL	0.890 <sub>0.001</sub>	0.890 <sub>0.001</sub>	0.890 <sub>0.001</sub>	0.876 <sub>0.001</sub>	0.840 <sub>0.004</sub>	0.880 <sub>0.004</sub>	0.877 <sub>0.001</sub>	0.893 <sub>0.000</sub>	0.890 <sub>0.001</sub>	0.896 <sub>0.001</sub>
	1-OE	0.646 <sub>0.009</sub>	0.647 <sub>0.010</sub>	0.646 <sub>0.008</sub>	0.533 <sub>0.008</sub>	0.511 <sub>0.016</sub>	0.639 <sub>0.009</sub>	0.609 <sub>0.007</sub>	0.667 <sub>0.006</sub>	0.654 <sub>0.007</sub>	0.683 <sub>0.006</sub>
	1-Cov	0.673 <sub>0.003</sub>	0.662 <sub>0.004</sub>	0.657 <sub>0.003</sub>	0.615 <sub>0.002</sub>	0.588 <sub>0.013</sub>	0.665 <sub>0.002</sub>	0.624 <sub>0.002</sub>	0.681 <sub>0.002</sub>	0.669 <sub>0.006</sub>	0.688 <sub>0.003</sub>
	1-RL	0.874 <sub>0.002</sub>	0.869 <sub>0.003</sub>	0.867 <sub>0.003</sub>	0.835 <sub>0.001</sub>	0.809 <sub>0.014</sub>	0.862 <sub>0.003</sub>	0.847 <sub>0.001</sub>	0.878 <sub>0.001</sub>	0.873 <sub>0.002</sub>	0.886 <sub>0.002</sub>
	AP	0.599 <sub>0.003</sub>	0.595 <sub>0.007</sub>	0.592 <sub>0.006</sub>	0.519 <sub>0.003</sub>	0.495 <sub>0.017</sub>	0.589 <sub>0.004</sub>	0.548 <sub>0.003</sub>	0.614 <sub>0.004</sub>	0.603 <sub>0.005</sub>	0.629 <sub>0.004</sub>
	AUC	0.861 <sub>0.001</sub>	0.855 <sub>0.002</sub>	0.854 <sub>0.002</sub>	0.828 <sub>0.001</sub>	0.801 <sub>0.017</sub>	0.852 <sub>0.003</sub>	0.839 <sub>0.001</sub>	0.864 <sub>0.001</sub>	0.859 <sub>0.002</sub>	0.871 <sub>0.002</sub>
	AVE	3.8	4.7	6.2	9.0	10.0	6.7	8.0	2.0	3.5	1.0
OBJ	1-HL	0.948 <sub>0.001</sub>	0.948 <sub>0.001</sub>	0.947 <sub>0.001</sub>	0.935 <sub>0.000</sub>	0.899 <sub>0.002</sub>	0.940 <sub>0.002</sub>	0.935 <sub>0.001</sub>	0.949 <sub>0.001</sub>	0.948 <sub>0.001</sub>	0.950 <sub>0.001</sub>
	1-OE	0.619 <sub>0.015</sub>	0.601 <sub>0.011</sub>	0.594 <sub>0.012</sub>	0.537 <sub>0.011</sub>	0.465 <sub>0.018</sub>	0.604 <sub>0.016</sub>	0.450 <sub>0.008</sub>	0.627 <sub>0.011</sub>	0.626 <sub>0.009</sub>	0.648 <sub>0.008</sub>
	1-Cov	0.807 <sub>0.006</sub>	0.794 <sub>0.006</sub>	0.793 <sub>0.006</sub>	0.768 <sub>0.005</sub>	0.744 <sub>0.008</sub>	0.796 <sub>0.008</sub>	0.759 <sub>0.006</sub>	0.813 <sub>0.006</sub>	0.809 <sub>0.006</sub>	0.818 <sub>0.006</sub>
	1-RL	0.888 <sub>0.005</sub>	0.876 <sub>0.004</sub>	0.875 <sub>0.004</sub>	0.860 <sub>0.004</sub>	0.833 <sub>0.006</sub>	0.878 <sub>0.006</sub>	0.850 <sub>0.004</sub>	0.890 <sub>0.005</sub>	0.889 <sub>0.004</sub>	0.897 <sub>0.003</sub>
	AP	0.639 <sub>0.010</sub>	0.627 <sub>0.009</sub>	0.623 <sub>0.010</sub>	0.577 <sub>0.008</sub>	0.512 <sub>0.014</sub>	0.630 <sub>0.012</sub>	0.537 <sub>0.008</sub>	0.649 <sub>0.009</sub>	0.649 <sub>0.008</sub>	0.665 <sub>0.006</sub>
	AUC	0.897 <sub>0.004</sub>	0.886 <sub>0.004</sub>	0.885 <sub>0.004</sub>	0.872 <sub>0.004</sub>	0.846 <sub>0.006</sub>	0.888 <sub>0.006</sub>	0.864 <sub>0.004</sub>	0.900 <sub>0.005</sub>	0.898 <sub>0.004</sub>	0.906 <sub>0.003</sub>
	AVE	4.0	5.7	6.8	8.2	9.8	5.3	9.0	2.0	3.0	1.0
PAS	1-HL	0.931 <sub>0.001</sub>	0.931 <sub>0.000</sub>	0.931 <sub>0.001</sub>	0.927 <sub>0.001</sub>	0.882 <sub>0.004</sub>	0.915 <sub>0.005</sub>	0.928 <sub>0.001</sub>	0.933 <sub>0.001</sub>	0.932 <sub>0.001</sub>	0.935 <sub>0.001</sub>
	1-OE	0.462 <sub>0.010</sub>	0.443 <sub>0.007</sub>	0.435 <sub>0.010</sub>	0.419 <sub>0.006</sub>	0.366 <sub>0.039</sub>	0.433 <sub>0.016</sub>	0.418 <sub>0.008</sub>	0.474 <sub>0.008</sub>	0.468 <sub>0.008</sub>	0.496 <sub>0.009</sub>
	1-Cov	0.781 <sub>0.007</sub>	0.749 <sub>0.003</sub>	0.738 <sub>0.010</sub>	0.720 <sub>0.004</sub>	0.674 <sub>0.011</sub>	0.759 <sub>0.006</sub>	0.738 <sub>0.003</sub>	0.790 <sub>0.006</sub>	0.778 <sub>0.004</sub>	0.795 <sub>0.005</sub>
	1-RL	0.830 <sub>0.006</sub>	0.804 <sub>0.002</sub>	0.792 <sub>0.008</sub>	0.778 <sub>0.003</sub>	0.736 <sub>0.011</sub>	0.808 <sub>0.006</sub>	0.797 <sub>0.002</sub>	0.836 <sub>0.006</sub>	0.828 <sub>0.004</sub>	0.844 <sub>0.004</sub>
	AP	0.549 <sub>0.007</sub>	0.517 <sub>0.004</sub>	0.510 <sub>0.008</sub>	0.482 <sub>0.005</sub>	0.438 <sub>0.022</sub>	0.524 <sub>0.009</sub>	0.486 <sub>0.005</sub>	0.562 <sub>0.005</sub>	0.552 <sub>0.006</sub>	0.581 <sub>0.007</sub>
	AUC	0.851 <sub>0.005</sub>	0.827 <sub>0.002</sub>	0.817 <sub>0.008</sub>	0.806 <sub>0.003</sub>	0.767 <sub>0.011</sub>	0.830 <sub>0.006</sub>	0.823 <sub>0.002</sub>	0.855 <sub>0.006</sub>	0.848 <sub>0.005</sub>	0.861 <sub>0.003</sub>
	AVE	3.5	5.7	7.2	8.7	10.0	6.0	7.5	2.0	3.5	1.0

### 3.5 ABLATION STUDY

Table 2: Ablation study on Pascal07, OBJECT and Mirflickr with PER=50% and LMR=50%. ‘✓’ and ‘✗’ represent the used and not used corresponding item, respectively.

$S_1$	$S_2$	$S_3$	Pascal07				OBJECT				Mirflickr			
			AP	AUC	1-RL	1-OE	AP	AUC	1-RL	1-OE	AP	AUC	1-RL	1-OE
✗	✓	✓	0.580	0.858	0.843	0.492	0.664	0.902	0.893	0.642	0.623	0.862	0.878	0.681
✓	✗	✓	0.561	0.845	0.834	0.476	0.652	0.896	0.889	0.632	0.612	0.851	0.865	0.665
✓	✓	✗	0.584	0.861	0.847	0.495	0.668	0.905	0.897	0.646	0.628	0.866	0.881	0.685
✓	✓	✓	<b>0.588</b>	<b>0.867</b>	<b>0.851</b>	<b>0.501</b>	<b>0.673</b>	<b>0.910</b>	<b>0.901</b>	<b>0.651</b>	<b>0.633</b>	<b>0.870</b>	<b>0.885</b>	<b>0.690</b>

The ablation studies are carried out to deeply examine the effect of the three key modules in TITRL, i.e., a dual-layer constraint framework for shared representation learning ( $S_1$ ), the strategy of view fusion guided by label information ( $S_2$ ), and sample-level label correlation semantic learning. After the separate removal of  $S_1$ ,  $S_2$ , and  $S_3$ , losses  $\mathcal{L}_f$  and  $\mathcal{L}_c$  tied to representation extraction are omitted, view aggregation is realized solely through distributed fusion, and loss  $\mathcal{L}_a$  is excluded without accounting for label dependencies, respectively. The ablation results shown in Table 2 lead



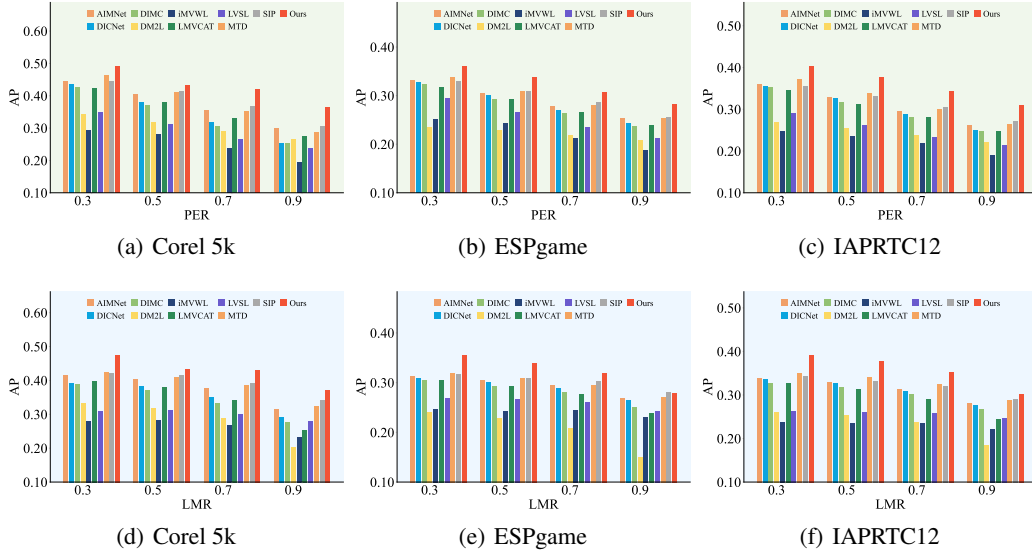


Figure 2: Experimental results on three datasets with one of PER and LMR fixed at 50% and the other varying from 30% to 90%

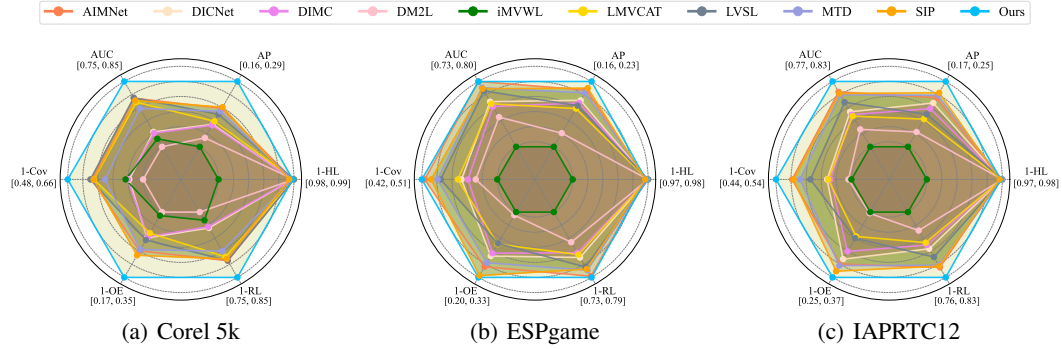


Figure 3: Experimental results of ten methods on three datasets with PER= 90% and LMR= 90%.

to the following findings: (i) The performance degradation observed upon the removal of any single module highlights the deliberate design of TITRL. (ii) Incorporating label semantics into the fusion process is instrumental in enhancing classification performance, as it enables the selective integration of the discriminative information from view representations. The training loss that facilitates the learning of both the shared representation and label correlation semantics exerts a positive influence, showing that our approach is valuable for advancing semantic exploration at both levels.

## 4 CONCLUSION

In this paper, we propose a Theory-Inspired Task-Relevant Representation Learning method called TITRL to address the IMvMLC problem, which is driven by the imperative to purify shared representations, improve the reliability of view fusion, and accurately capture multi-label correlation semantics. Specifically, TITRL introduces a dual-layer constraint framework based on the interaction of mutual information, which enables the disentanglement of view-specific noise. By deriving the tractable variational bounds, we provide theoretical guidance for learning pure shared information in a principled manner. For view fusion, TITRL integrates a distribution-aware strategy that leverages the statistical view property with a confidence-driven late fusion mechanism, thereby reinforcing the stability of representation expression and predictive reliability. Furthermore, we explicitly model sample-level label correlations by aligning shared representations with learnable label prototypes, allowing for flexible use of label dependencies. Extensive experiments on multiple benchmark datasets demonstrate TITRL’s superiority and robustness, particularly in high-deficiency conditions where most baselines fail.

## ETHICS STATEMENT

We confirm that our work adheres to the ICLR Code of Ethics. This research does not involve human subjects or raise concerns regarding privacy, legal compliance, or harmful insights. The datasets used are publicly available and properly cited. No conflicts of interest or external sponsorships influenced the research. We are committed to maintaining research integrity throughout the process.

## REPRODUCIBILITY STATEMENT

We have taken several steps to ensure the reproducibility of our work. The detailed description of the methodology, experimental setup, and datasets used can be found in the main text and supplementary materials. The key components required for reproducing the experiments include the model architecture, training procedure, and evaluation metrics, all of which are explicitly described. In addition, for the reproducibility of the novel models and algorithms, we will release the source code upon acceptance. For the theoretical claims, we provide a comprehensive explanation of the underlying assumptions and the complete proofs of the claims in the appendix. Furthermore, all datasets used in the experiments are publicly available, and a detailed description of the data processing steps, including handling missing views and labels, can be found in the supplementary materials. We encourage the readers to refer to the provided supplementary materials for any further details required to replicate the results.

## REFERENCES

- Luis Von Ahn and Laura Dabbish. Labeling images with a computer game. In *SIGCHI Conference on Human Factors in Computing Systems*, pp. 319–326, 2004.
- Wei Chen, Jiage Chen, Yuewu Wan, Xining Liu, Mengya Cai, Jingguo Xu, Hongbo Cui, and Mengdie Duan. Land cover classification based on multimodal remote sensing fusion. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, 10:35–40, 2024.
- Pinar Duygulu, Kobus Barnard, Joao FG de Freitas, and David A Forsyth. Object recognition as machine translation: Learning a lexicon for a fixed image vocabulary. In *Proceedings of European Conference on Computer Vision*, pp. 97–112, 2002.
- Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88: 303–338, 2010.
- Xiang Fang, Yuchong Hu, Pan Zhou, and Dapeng Wu. Animc: A soft approach for autoweighted noisy and incomplete multiview clustering. *IEEE Transactions on Artificial Intelligence*, 3(2): 192–206, 2021.
- Xiang Fang, Daizong Liu, Pan Zhou, and Yuchong Hu. Multi-modal cross-domain alignment network for video moment retrieval. *IEEE Transactions on Multimedia*, 25:7517–7532, 2022.
- Marco Federici, Anjan Dutta, Patrick Forré, Nate Kushman, and Zeynep Akata. Learning robust representations via multi-view information bottleneck. In *ICLR*, 2020.
- Michael Grubinger, Paul Clough, Henning Müller, and Thomas Deselaers. The iapr tc-12 benchmark: A new evaluation resource for visual information systems. In *International Workshop onto Image*, volume 2, pp. 1–11, 2006.
- Renxiang Guan, Tianrui Liu, Wenxuan Tu, Chang Tang, Wenhan Luo, and Xinwang Liu. Sampling enhanced contrastive multi-view remote sensing data clustering with long-short range information mining. *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–15, 2025.
- Pingting Hao, Kunpeng Liu, and Wanfu Gao. Anchor-guided global view reconstruction for multi-view multi-label feature selection. *Information Sciences*, 679:121124, 2024.
- Pingting Hao, Wanfu Gao, and Liang Hu. Embedded feature fusion for multi-view multi-label feature selection. *Pattern Recognition*, 157:110888, 2025.

- Mark J Huiskes and Michael S Lew. The mir flickr retrieval evaluation. In *Proceedings of ACM International Conference on Multimedia Information Retrieval*, pp. 39–43, 2008.
- Yinggang Li, Xiangrong Tong, and Zhongming Lv. Multi-dimensional requirements for reinforcement recommendation reasoning. *Applied Intelligence*, 55(6):1–16, 2025.
- Chengliang Liu, Zhihao Wu, Jie Wen, Yong Xu, and Chao Huang. Localized sparse incomplete multi-view clustering. *IEEE Transactions on Multimedia*, 25:5539–5551, 2022.
- Chengliang Liu, Jie Wen, Yabo Liu, Chao Huang, Zhihao Wu, Xiaoling Luo, and Yong Xu. Masked two-channel decoupling framework for incomplete multi-view weak multi-label learning. *Advances in Neural Information Processing Systems*, 36:32387–32400, 2023a.
- Chengliang Liu, Jie Wen, Xiaoling Luo, Chao Huang, Zhihao Wu, and Yong Xu. Dicnet: Deep instance-level contrastive network for double incomplete multi-view multi-label classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 8807–8815, 2023b.
- Chengliang Liu, Jie Wen, Xiaoling Luo, and Yong Xu. Incomplete multi-view multi-label learning via label-guided masked view-and category-aware transformers. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 8816–8824, 2023c.
- Chengliang Liu, Jinlong Jia, Jie Wen, Yabo Liu, Xiaoling Luo, Chao Huang, and Yong Xu. Attention-induced embedding imputation for incomplete multi-view partial multi-label classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 13864–13872, 2024a.
- Chengliang Liu, Gehui Xu, Jie Wen, Yabo Liu, Chao Huang, and Yong Xu. Partial multi-view multi-label classification via semantic invariance learning and prototype modeling. In *Forty-first International Conference on Machine Learning*, 2024b.
- Chengliang Liu, Jie Wen, Yong Xu, Bob Zhang, Liqiang Nie, and Min Zhang. Reliable representation learning for incomplete multi-view missing multi-label classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–17, 2025.
- Meng Liu, Yong Luo, Dacheng Tao, Chao Xu, and Yonggang Wen. Low-rank multi-view learning in matrix completion for multi-label image classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 29, 2015.
- Yuanyuan Liu, Jiyao Peng, Wei Dai, Jiabei Zeng, and Shiguang Shan. Joint spatial and scale attention network for multi-view facial expression recognition. *Pattern Recognition*, 139:109496, 2023d.
- Zhongchen Ma and Songcan Chen. Expand globally, shrink locally: Discriminant multi-label learning with missing labels. *Pattern Recognition*, 111:107675, 2021.
- Yalan Qin, Xinpeng Zhang, Shui Yu, and Guorui Feng. A survey on representation learning for multi-view data. *Neural Networks*, 181:106842, 2025.
- Chongjie Si, Yuheng Jia, Ran Wang, Min-Ling Zhang, Yanghe Feng, and Chongxiao Qu. Multi-label classification with high-rank and high-order label correlations. *IEEE Transactions on Knowledge and Data Engineering*, 36(8):4076–4088, 2023.
- Gencer Sumbul, Marcela Charfuelan, Begüm Demir, and Volker Markl. Bigearthnet: A large-scale benchmark archive for remote sensing image understanding. In *IGARSS 2019-2019 IEEE international geoscience and remote sensing symposium*, pp. 5901–5904. IEEE, 2019.
- Zhaobin Sun, Nannan Wu, Junjie Shi, Li Yu, Kwang-Ting Cheng, and Zengqiang Yan. Fedmlp: Federated multi-label medical image classification under task heterogeneity. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 394–404. Springer, 2024.
- Qiaoyu Tan, Guoxian Yu, Carlotta Domeniconi, Jun Wang, and Zili Zhang. Incomplete multi-view weak-label learning. In *Ijcai*, pp. 2703–2709, 2018.

- Zhihui Tian, John Upchurch, G Austin Simon, José Dubeux, Alina Zare, Chang Zhao, and Joel B Harley. Quantifying heterogeneous ecosystem services with multi-label soft classification. In *IGARSS 2024-2024 IEEE International Geoscience and Remote Sensing Symposium*, pp. 427–431. IEEE, 2024.
- Jie Wen, Zheng Zhang, Yong Xu, Bob Zhang, Lunke Fei, and Hong Liu. Unified embedding alignment with missing views inferring for incomplete multi-view clustering. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pp. 5393–5400, 2019.
- Jie Wen, Zheng Zhang, Lunke Fei, Bob Zhang, Yong Xu, Zhao Zhang, and Jinxing Li. A survey on incomplete multiview clustering. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 53(2):1136–1149, 2022.
- Jie Wen, Chengliang Liu, Shijie Deng, Yicheng Liu, Lunke Fei, Ke Yan, and Yong Xu. Deep double incomplete multi-view multi-label learning with incomplete labels and missing views. *IEEE transactions on neural networks and learning systems*, 2023.
- Xuan Wu, Qing-Guo Chen, Yao Hu, Dengbao Wang, Xiaodong Chang, Xiaobo Wang, and Min-Ling Zhang. Multi-view multi-label learning with view-specific information extraction. In *IJCAI*, pp. 3884–3890, 2019.
- Changqing Zhang, Ziwei Yu, Qinghua Hu, Pengfei Zhu, Xinwang Liu, and Xiaobo Wang. Latent semantic aware multi-view multi-label classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Dawei Zhao, Qingwei Gao, Yixiang Lu, Dong Sun, and Yusheng Cheng. Consistency and diversity neural network multi-view multi-label learning. *Knowledge-Based Systems*, 218:106841, 2021.
- Dawei Zhao, Qingwei Gao, Yixiang Lu, and Dong Sun. Non-aligned multi-view multi-label classification via learning view-specific labels. *IEEE Transactions on Multimedia*, 25:7235–7247, 2022.
- Lihua Zhou, Guowang Du, Kevin Lue, Lizheng Wang, and Jingwei Du. A survey and an empirical evaluation of multi-view clustering approaches. *ACM Computing Surveys*, 56(7):1–38, 2024.
- Yue Zhu, James T Kwok, and Zhi-Hua Zhou. Multi-label learning with global and local label correlation. *IEEE Transactions on Knowledge and Data Engineering*, 30(6):1081–1094, 2017.

## A APPENDIX

### A.1 RELATED WORKS

#### A.1.1 MvMLC

The integration of multi-view learning with multi-label classification introduces algorithmic complexity but enables a comprehensive representation of object attributes. Consequently, algorithms in this domain were designed to address multi-view processing and multi-label recognition concurrently. For instance, LSA-MML Zhang et al. (2018) developed a shared latent representation via matrix factorization and aligns latent basis matrices across views in the kernel space using the Hilbert-Schmidt Independence Criterion (HSIC) to maximize inter-view dependency. Wu et al. proposed SIMM Wu et al. (2019), a deep MvMLC network that learned a shared subspace and view-specific features, ensuring a clear separation between shared and specific representations to capture the distinct contribution of each view. CDMM Zhao et al. (2021) focused on label relevance by promoting inter-view consistency and diversity through a straightforward approach, constructing a label affinity matrix that integrates Jaccard similarity in the label space with instance-level proximity in the feature space, followed by label propagation to enhance label representation. LVSL Zhao et al. (2022) addressed non-aligned multi-view multi-label classification by jointly exploring view-specific labels and low-rank label structures while preserving geometric properties of the original data via Laplacian graph regularization.

#### A.1.2 iMvMLC

Incomplete multi-view multi-label learning has attracted growing attention due to challenges posed by missing data in both feature and label domains. Missing views are generally addressed either by disregarding them with prior knowledge or reconstructing them through data completion. Liu et al. (2022) introduced a missing indicator matrix that enabled adaptively handling instances with missing views, while Wen et al. (2019) proposed UEAF, which reconstructed missing views using an error matrix informed by local structure and enhanced alignment through reverse graph learning. The challenge of incomplete labels has inspired multiple solutions. DM2L Ma & Chen (2021) inferred latent labels from partial annotations, integrated label correlations and low-rank constraints for robust feature learning, and refined pseudo-labels via self-supervised learning. GLOCAL Zhu et al. (2017) constructed bidirectional graphs between labels and instances to capture global label semantics under missing annotations, incorporated local priors and structure-preserving constraints to mitigate label noise and enhance feature discriminability. The methods capable of simultaneously handling missing views and labels initially involved the exploration of some traditional approaches. For example, iMVWL Tan et al. (2018) handled missing views and labels by jointly learning a shared subspace, a weak-label predictor, and a low-rank label correlation matrix. NAIM3L captured global high-rank and local low-rank structures of multi-label matrix and employed an efficient ADMM algorithm with linear time complexity for large-scale data. However, the performance of these methods is limited by their capacity to extract shallow representations and perform linear predictions.

Recent deep learning-based methods have achieved substantial improvements. DIMC Wen et al. (2023) utilizes an end-to-end network to extract high-level discriminative features with decoder-based reconstruction for robustness. DICNet Liu et al. (2023b) applied instance-level contrastive learning to unify representations of identical samples across views. LMVCAT Liu et al. (2023c) employed a Transformer-based framework with multi-head attention to facilitate feature interactions and category-aware modules to capture label semantics. Additional methods, including AIMNet Liu et al. (2024a), SIP Liu et al. (2024b), and MTD Liu et al. (2023a), explored attention-guided embedding completion, information bottleneck principles, and masked dual-channel decoupling strategies to reconstruct missing views, derive shared information and capture view-specific information, respectively. Together, these approaches emphasize the significance of extracting advanced representations, preserving label correlations, and utilizing sophisticated network architectures to address incomplete multi-view multi-label learning challenges.

## A.2 A COMPLETE DERIVATION OF THE TASK-RELEVANT REPRESENTATION LEARNING MODEL UNDER A DUAL-LAYER CONSTRAINT

In this section, we present a detailed derivation of model (4):

$$\min \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \left( \underbrace{-I(\mathbf{x}^{(v)}; \mathbf{z}) + I(\{\mathbf{x}^{\sim(v)}\}; \mathbf{z} | \mathbf{x}^{(v)})}_{\text{feature-level}} + \underbrace{I(\mathbf{x}^{(v)}; \mathbf{y}) - I(\mathbf{z}^{(v)}; \mathbf{y}) - I(\mathbf{y}; \mathbf{z}^{(v)}) + I(\mathbf{y}; \mathbf{z}^{(v)} | \{\mathbf{z}^{\sim(v)}\})}_{\text{category-level}} \right). \quad (21)$$

For the first term, we proceed with the following expansion based on the definition of mutual information:

$$I(\mathbf{x}^{(v)}; \mathbf{z}) = \int \int p(\mathbf{x}^{(v)}, \mathbf{z}) \log \frac{p(\mathbf{x}^{(v)}, \mathbf{z})}{p(\mathbf{x}^{(v)})p(\mathbf{z})} d\mathbf{z} d\mathbf{x}^{(v)}. \quad (22)$$

Considering  $p(\mathbf{x}^{(v)}, \mathbf{z}) = p(\mathbf{x}^{(v)} | \mathbf{z})p(\mathbf{x}^{(v)})$ , we have

$$\begin{aligned} I(\mathbf{x}^{(v)}; \mathbf{z}) &= \int \int p(\mathbf{x}^{(v)} | \mathbf{z})p(\mathbf{x}^{(v)}) \log \frac{p(\mathbf{x}^{(v)} | \mathbf{z})p(\mathbf{x}^{(v)})}{p(\mathbf{x}^{(v)})p(\mathbf{z} | \mathbf{x}^{(v)})} d\mathbf{z} d\mathbf{x}^{(v)} \\ &= \int \int p(\mathbf{x}^{(v)} | \mathbf{z})p(\mathbf{x}^{(v)}) \log p(\mathbf{x}^{(v)} | \mathbf{z}) d\mathbf{z} d\mathbf{x}^{(v)} \\ &\quad + \int \int p(\mathbf{x}^{(v)} | \mathbf{z})p(\mathbf{z}) \log \frac{1}{p(\mathbf{z})} d\mathbf{z} d\mathbf{x}^{(v)}. \end{aligned} \quad (23)$$

Since  $H(\mathbf{z}) = -\int p(\mathbf{z}) \log p(\mathbf{z}) d\mathbf{z} \geq 0$ , we have

$$\begin{aligned} I(\mathbf{x}^{(v)}; \mathbf{z}) &= \int \int p(\mathbf{x}^{(v)} | \mathbf{z})p(\mathbf{x}^{(v)}) \log p(\mathbf{x}^{(v)} | \mathbf{z}) d\mathbf{z} d\mathbf{x}^{(v)} \\ &\quad + \int p(\mathbf{x}^{(v)} | \mathbf{z}) H(\mathbf{z}) d\mathbf{z} \\ &\geq \int \int p(\mathbf{x}^{(v)} | \mathbf{z})p(\mathbf{x}^{(v)}) \log p(\mathbf{x}^{(v)} | \mathbf{z}) d\mathbf{z} d\mathbf{x}^{(v)} \\ &= \int \int p(\mathbf{x}^{(v)} | \mathbf{z})p(\mathbf{x}^{(v)}) \log q^v(\mathbf{x}^{(v)} | \mathbf{z}) d\mathbf{z} d\mathbf{x}^{(v)} \\ &\quad + \int \int p(\mathbf{x}^{(v)} | \mathbf{z})p(\mathbf{x}^{(v)}) \log \frac{p(\mathbf{x}^{(v)} | \mathbf{z})}{q^v(\mathbf{x}^{(v)} | \mathbf{z})} d\mathbf{z} d\mathbf{x}^{(v)}. \end{aligned} \quad (24)$$

Based on the definition of the Kullback-Leibler divergence, we can get

$$D_{KL}(p(\mathbf{x}^{(v)} | \mathbf{z}) \| q^v(\mathbf{x}^{(v)} | \mathbf{z})) = \int p(\mathbf{x}^{(v)} | \mathbf{z}) \log \frac{p(\mathbf{x}^{(v)} | \mathbf{z})}{q^v(\mathbf{x}^{(v)} | \mathbf{z})} d\mathbf{x}^{(v)}. \quad (25)$$

Since  $D_{KL}(p(\mathbf{x}^{(v)} | \mathbf{z}) \| q^v(\mathbf{x}^{(v)} | \mathbf{z})) \geq 0$ , we have

$$\begin{aligned} I(\mathbf{x}^{(v)}; \mathbf{z}) &\geq \int \int p(\mathbf{x}^{(v)} | \mathbf{z})p(\mathbf{x}^{(v)}) \log q^v(\mathbf{x}^{(v)} | \mathbf{z}) d\mathbf{z} d\mathbf{x}^{(v)} \\ &\quad + \int p(\mathbf{x}^{(v)}) D_{KL}(p(\mathbf{x}^{(v)} | \mathbf{z}) \| q^v(\mathbf{x}^{(v)} | \mathbf{z})) d\mathbf{z} \\ &\geq \int \int p(\mathbf{x}^{(v)} | \mathbf{z})p(\mathbf{x}^{(v)}) \log q^v(\mathbf{x}^{(v)} | \mathbf{z}) d\mathbf{z} d\mathbf{x}^{(v)} \\ &= \int \int p(\mathbf{z}, \mathbf{x}^{(v)}) \log q^v(\mathbf{x}^{(v)} | \mathbf{z}) d\mathbf{z} d\mathbf{x}^{(v)}. \end{aligned} \quad (26)$$

According to the property of integration, we have  $\int \int p(z, \mathbf{x}^{(v)}) \log q^v(\mathbf{x}^{(v)}|z) dz d\mathbf{x}^{(v)} = \int \int p(\{\mathbf{x}\}, z) \log q^v(\mathbf{x}^{(v)}|z) d\{\mathbf{x}\} dz$ . Thus, Eq. (26) can be rewritten as

$$\begin{aligned} & I(\mathbf{x}^{(v)}; z) \\ & \geq \int \int p(\{\mathbf{x}\}, z) \log q^v(\mathbf{x}^{(v)}|z) d\{\mathbf{x}\} dz \\ & = \int p(\{\mathbf{x}\}) \int p(z|\{\mathbf{x}\}) \log q^v(\mathbf{x}^{(v)}|z) d\{\mathbf{x}\} dz \\ & = \mathbb{E}_{\mathbf{x} \sim p(\{\mathbf{x}\})} \left[ \int p(z|\{\mathbf{x}\}) \log q^v(\mathbf{x}^{(v)}|z) dz \right] \end{aligned} \quad (27)$$

Regarding the second term  $I(\{\mathbf{x}^{(v)}\}; z | \mathbf{x}^{(v)})$ , we have the following expansion based on the definition of conditional mutual information  $I(a; b | c) = \iiint p(a, b, c) \log \left( \frac{p(a, b|c)}{p(a|c)p(b|c)} \right) dadbdc$ :

$$\begin{aligned} & I(\{\mathbf{x}^{(v)}\}; z | \mathbf{x}^{(v)}) \\ & = \int \int p(\{\mathbf{x}\}, z) \log \frac{p(\{\mathbf{x}\}, z) p(\mathbf{x}^{(v)})}{p(\{\mathbf{x}\}) p(z, \mathbf{x}^{(v)})} d\{\mathbf{x}\} dz \\ & = \int \int p(\{\mathbf{x}\}, z) \log \frac{p(z|\{\mathbf{x}\})}{p(z|\mathbf{x}^{(v)})} d\{\mathbf{x}\} dz \end{aligned} \quad (28)$$

By introducing the variational distribution  $g^v(z|\mathbf{x}^{(v)})$  to estimate  $p(z|\mathbf{x}^{(v)})$ , we have

$$\begin{aligned} & I(\{\mathbf{x}^{(v)}\}; z | \mathbf{x}^{(v)}) \\ & = \int \int p(\{\mathbf{x}\}, z) \log \frac{p(z|\{\mathbf{x}\})}{g^v(z|\mathbf{x}^{(v)})} d\{\mathbf{x}\} dz + \int \int p(\{\mathbf{x}\}, z) \log \frac{g^v(z|\mathbf{x}^{(v)})}{p(z|\mathbf{x}^{(v)})} d\{\mathbf{x}\} dz \\ & = \int \int p(\{\mathbf{x}\}, z) \log \frac{p(z|\{\mathbf{x}\})}{g^v(z|\mathbf{x}^{(v)})} d\{\mathbf{x}\} dz + \int p(\{\mathbf{x}\}) \int p(z|\{\mathbf{x}\}) \log \frac{g^v(z|\mathbf{x}^{(v)})}{p(z|\mathbf{x}^{(v)})} d\{\mathbf{x}\} dz \\ & = \int \int p(\{\mathbf{x}\}, z) \log \frac{p(z|\{\mathbf{x}\})}{g^v(z|\mathbf{x}^{(v)})} d\{\mathbf{x}\} dz + \int p(\mathbf{x}^{(v)}) \int p(z|\mathbf{x}^{(v)}) \log \frac{g^v(z|\mathbf{x}^{(v)})}{p(z|\mathbf{x}^{(v)})} d\mathbf{x}^{(v)} dz \\ & = \int \int p(\{\mathbf{x}\}, z) \log \frac{p(z|\{\mathbf{x}\})}{g^v(z|\mathbf{x}^{(v)})} d\{\mathbf{x}\} dz - \int p(\mathbf{x}^{(v)}) D_{KL}(p(z|\mathbf{x}^{(v)}) || g^v(z|\mathbf{x}^{(v)})) d\mathbf{x}^{(v)} \\ & \leq \int \int p(\{\mathbf{x}\}, z) \log \frac{p(z|\{\mathbf{x}\})}{g^v(z|\mathbf{x}^{(v)})} d\{\mathbf{x}\} dz \\ & = \mathbb{E}_{\{\mathbf{x}\} \sim p(\{\mathbf{x}\})} \left[ D_{KL} \left( p(z|\{\mathbf{x}\}) || g^v(z|\mathbf{x}^{(v)}) \right) \right], \end{aligned} \quad (29)$$

where  $p(z|\{\mathbf{x}\})$  is determined by the distribution fusion method provided in Eq. 14

Regarding the category-level constraint  $\min I(\mathbf{x}^{(v)}; \mathbf{y}) - I(\mathbf{z}^{(v)}; \mathbf{y})$  in the of model (4), we have the following equivalent transformation according to the equality  $I(\mathbf{x}^{(v)}; \mathbf{y}) = H(\mathbf{y}) - H(\mathbf{y}|\mathbf{x}^{(v)})$

$$\min I(\mathbf{x}^{(v)}; \mathbf{y}) - I(\mathbf{z}^{(v)}; \mathbf{y}) \iff \min H(\mathbf{y}|\mathbf{z}^{(v)}) - H(\mathbf{y}|\mathbf{x}^{(v)}). \quad (30)$$

Utilizing the definition of entropy, we can get

$$H(\mathbf{y} | \mathbf{z}^{(v)}) = - \int p(\mathbf{y} | \mathbf{z}^{(v)}) \log p(\mathbf{y} | \mathbf{z}^{(v)}) d\mathbf{y}. \quad (31)$$

Thus, we can find that  $H(\mathbf{y} | \mathbf{z}^{(v)})$  depends solely on  $p(\mathbf{y} | \mathbf{z}^{(v)})$ . As a result, minimizing  $I(\mathbf{x}^{(v)}; \mathbf{y}) - I(\mathbf{z}^{(v)}; \mathbf{y})$  naturally transforms into optimizing the distance between the corresponding two distributions:

$$\min \sum_{v \in \mathcal{V}} D_{KL} \left( p(\mathbf{y}|\mathbf{z}^{(v)}) || p(\mathbf{y}|\mathbf{x}^{(v)}) \right). \quad (32)$$

For the last term  $I(\mathbf{y}; \mathbf{z}^{(v)} | \{\mathbf{z}^{\sim(v)}\})$ , we have

$$\begin{aligned}
& I(\mathbf{y}; \mathbf{z}_i | \{\mathbf{z}^{\sim(v)}\}) \\
&= H(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) - H(\mathbf{y} | \{\mathbf{z}\}) \\
&= - \int p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \log p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) d\mathbf{y} + \int p(\mathbf{y} | \{\mathbf{z}\}) \log p(\mathbf{y} | \{\mathbf{z}\}) d\mathbf{y} \\
&= - \int p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \log \left[ \frac{p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})}{p(\mathbf{y} | \{\mathbf{z}\})} p(\mathbf{y} | \{\mathbf{z}\}) \right] d\mathbf{y} \\
&\quad + \int p(\mathbf{y} | \{\mathbf{z}\}) \log \left[ \frac{p(\mathbf{y} | \{\mathbf{z}\})}{p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})} p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \right] d\mathbf{y},
\end{aligned} \tag{33}$$

By adding terms in the logarithmic operation, we can obtain:

$$\begin{aligned}
& I(\mathbf{y}; \mathbf{z}_i | \{\mathbf{z}^{\sim(v)}\}) \\
&= - \int p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \log \left[ \frac{p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})}{p(\mathbf{y} | \{\mathbf{z}\})} p(\mathbf{y} | \{\mathbf{z}\}) \right] d\mathbf{y} \\
&\quad + \int p(\mathbf{y} | \{\mathbf{z}\}) \log \left[ \frac{p(\mathbf{y} | \{\mathbf{z}\})}{p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})} p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \right] d\mathbf{y} \\
&= - \int p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \log \left[ \frac{p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})}{p(\mathbf{y} | \{\mathbf{z}\})} \right] d\mathbf{y} - \int p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \log p(\mathbf{y} | \{\mathbf{z}\}) d\mathbf{y} \\
&\quad + \int p(\mathbf{y} | \{\mathbf{z}\}) \log \left[ \frac{p(\mathbf{y} | \{\mathbf{z}\})}{p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})} \right] d\mathbf{y} + \int p(\mathbf{y} | \{\mathbf{z}\}) \log p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) d\mathbf{y}. \\
&= - D_{KL}(p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) || p(\mathbf{y} | \{\mathbf{z}\})) + H(p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}); p(\mathbf{y} | \{\mathbf{z}\})) \\
&\quad + D_{KL}(p(\mathbf{y} | \{\mathbf{z}\}) || p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})) - H(p(\mathbf{y} | \{\mathbf{z}\}); p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})) \\
&\leq D_{KL}(p(\mathbf{y} | \{\mathbf{z}\}) || p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})) + H(p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}); p(\mathbf{y} | \{\mathbf{z}\})).
\end{aligned} \tag{34}$$

Since the optimization goals of  $D_{KL}(p(\mathbf{y} | \{\mathbf{z}\}) || p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}))$  and  $H(p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}); p(\mathbf{y} | \{\mathbf{z}\}))$  are both to make the distributions of  $p(\mathbf{y} | \{\mathbf{z}\})$  and  $p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})$  closer, we adopt the tractable term  $D_{KL}(p(\mathbf{y} | \{\mathbf{z}\}) || p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}))$  to minimize  $I(\mathbf{y}; \mathbf{z}^{(v)} | \{\mathbf{z}^{\sim(v)}\})$ . For the term  $I(\mathbf{y}; \mathbf{z}^{(v)})$  in model (4), its purpose is to enhance the information correlation between the representation and the labels, which can be achieved by optimizing the classification loss. In conclusion, the final derived training loss is as follows:

$$\begin{cases} \mathcal{L}_f = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \left[ -\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z} | \{\mathbf{x}\})} \log q^v(\mathbf{x}^{(v)} | \mathbf{z}) + D_{KL}(p(\mathbf{z} | \{\mathbf{x}\}) || g^v(\mathbf{z} | \mathbf{x}^{(v)})) \right] \\ \mathcal{L}_c = \frac{1}{|\mathcal{V}|} \sum_{v \in \mathcal{V}} \left( D_{KL}(p(\mathbf{y} | \mathbf{z}^{(v)}) || p(\mathbf{y} | \mathbf{x}^{(v)})) + D_{KL}(p(\mathbf{y} | \{\mathbf{z}\}) || p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})) \right). \end{cases} \tag{35}$$

### A.3 TRAINING PROCESS

The main training process of SMVTEP is summarised in Algorithm 1.

### A.4 EXPERIMENT

#### A.4.1 EXPERIMENT SETUP

**Datasets and Comparison Methods.** We conduct experiments on six publicly available multi-view multi-label datasets, summarized in Table 3. The details of these datasets are as follows. **Corel 5k** contains 4999 images with 260 annotation words, where each word can be treated as a semantic label. **IAPRTC12** consists of 19,627 high-resolution natural images with 261 possible labels, covering categories such as sports, animals, human activities, and urban scenes. **ESPGAME** provides 20,770 images with 268 associated tags collected from an online image labeling game. **Pascal07**



**Algorithm 1** Training process of TITRL

**Input:** Incomplete multi-view multi-label data  $(\{\mathbf{x}^{(v)}\}_{v=1}^V, \mathbf{y})$ , observed view set  $\mathcal{V}$ , known label set  $\mathcal{U}$ , and balanced parameters  $\lambda_1$  and  $\lambda_2$ .

**Output:** Prediction  $\mathbf{p}^t$ .

- 1: Initialize the parameters  $\{\mathbf{b}_i\}_{i=1}^c$  as  $\mathbf{I} \in \{0, 1\}^{c \times c}$
- 2: **for**  $t = 0$ ;  $t < \text{Total epoch}$ ;  $t++$  **do**
- 3:   Compute the variational distribution  $\{g^v(\mathbf{z} | \mathbf{x}^{(v)}) | v \in \mathcal{V}\}$  of the shared representations extracted from each available view through encoders  $f_\mu^v$  and  $f_{\sigma^2}^v$ .
- 4:   Compute the distribution parameters of the integrated shared representation  $\mathbf{z}$  by Eq. (14)
- 5:   Sample  $\mathbf{z}$  from the corresponding distribution by Eq. (15).
- 6:   Compute the conditional distribution  $\{q^v(\mathbf{x}^{(v)} | \mathbf{z}) | v \in \mathcal{V}\}$  by decoders
- 7:   Obtain the distribution of  $c$  label prototype  $\{\mathbf{l}_i | \mathbf{l}_i \sim \mathcal{N}(h_\mu(\mathbf{b}_i), h_{\sigma^2}(\mathbf{b}_i)\mathbf{I}) | i = 1, \dots, c\}$ .
- 8:   Sample each  $\mathbf{l}_i$  from the corresponding distribution like Eq. (15).
- 9:   Compute the pseudo-predictions  $p(\mathbf{y}_i | \{\mathbf{z}^{(v)}\})$ ,  $p(\mathbf{y}_i | \mathbf{x}^{(v)})$ , and  $p(\mathbf{y}_i | \mathbf{z}^{(v)})$  by Eq. (15).
- 10:   Compute the final prediction  $\mathbf{p}^t$  by Eq. (18)
- 11:   Compute the total loss  $\mathcal{L}$  by Eq. (20).
- 12:   Update network parameters.
- 13: **end for**

is a widely used benchmark for object detection and recognition, comprising 9963 images across 20 object categories. **Mirflickr** includes 25,000 Flickr images, each annotated with up to 38 labels. **OBJECT** contains 6047 samples described from five different perspectives and annotated with 31 attributes. To demonstrate the effectiveness of our method, we compare it with nine representative approaches: AIMNet Liu et al. (2024a), DICNet Liu et al. (2023b), DIMC Wen et al. (2023), iMVWL Tan et al. (2018), LMVCAT Liu et al. (2023c), MTD Liu et al. (2023a), SIP Liu et al. (2024b), LVSL Zhao et al. (2022), and DM2L Ma & Chen (2021). A detailed summary of these comparison methods, including their sources and functionality, is presented in Table 4.

Table 3: Detailed information of datasets.

View	Object	VOC 2007	Corel 5k	ESP Game	IAPR TC-12	MIR Flickr
1	CH (64)	DenseHue (100)	DenseHue (100)	DenseHue (100)	DenseHue (100)	DenseHue (100)
2	CM (225)	DenseSift (1000)	DenseSift (1000)	DenseSift (1000)	DenseSift (1000)	DenseSift (1000)
3	CORR (144)	GIST (512)	GIST (512)	GIST (512)	GIST (512)	GIST (512)
4	EDH (73)	HSV (4096)	HSV (4096)	HSV (4096)	HSV (4096)	HSV (4096)
5	WT (128)	RGB (4096)	RGB (4096)	RGB (4096)	RGB (4096)	RGB (4096)
6	—	LAB (4096)	LAB (4096)	LAB (4096)	LAB (4096)	LAB (4096)
#Labels	31	20	260	268	291	38
#Instances	6047	9963	4999	20770	19627	25000

Table 4: Detailed information of comparison methods. ✓ indicates the method can handle the corresponding problem, while ✗ denotes it cannot.

Method	Source	Year	Multi-label	Multi-view	Missing-view	Missing-label
iMVWL	IJCAI	2018	✓	✓	✓	✓
DM2L	PR	2021	✓	✗	✗	✓
LVLS	TMM	2022	✓	✓	✗	✗
DICNet	AAAI	2023	✓	✓	✓	✓
DIMC	TNNLS	2023	✓	✓	✓	✓
LMVCAT	AAAI	2023	✓	✓	✓	✓
AIMNet	AAAI	2024	✓	✓	✓	✓
MTD	NeurIPS	2024	✓	✓	✓	✓
SIP	ICML	2024	✓	✓	✓	✓

**Implementation Details.** We adopt six widely used evaluation metrics to ensure consistency across experiments: Hamming Loss (HL), Ranking Loss (RL), OneError (OE), Coverage (Cov), Average

Precision (AP), and Area Under Curve (AUC). In general, higher values of AP and AUC indicate stronger predictive ability, while lower HL, RL, OE, and Cov reflect better classification accuracy. Specifically: (1) HL measures the proportion of incorrectly predicted labels; (2) RL quantifies the proportion of label pairs that are incorrectly ordered; (3) OE checks whether the top-ranked predicted label matches the ground truth; (4) Cov evaluates how many predicted labels need to be considered to cover all true labels; (5) AP corresponds to the average precision across different recall levels; and (6) AUC represents the probability that a randomly chosen positive instance is ranked above a randomly chosen negative one. For all datasets, the neighbor number  $k$  is set to 10. Optimization is performed using the Adam optimizer with an initial learning rate of 0.0001. All models share the same dataset partitioning, and the positions of missing views and labels are fixed across methods to guarantee fairness in comparison.

#### A.4.2 EXPERIMENT RESULT

##### Comparative Experiment.

We provide the complete results on six datasets by varying PER and LMR from 30% to 90%. Fig. 4 reports the case where PER changes while LMR is fixed at 50%. As expected, performance gradually decreases for all methods as more views are removed, but TITRL consistently secures the best accuracy. The performance gap over baselines becomes more significant at higher missing levels, indicating that TITRL is particularly effective at mitigating the negative impact of severe view absence through robust multi-view representation learning. Fig. 5 presents the results of varying LMR with PER fixed at 50%. Compared with the PER case, the decline under missing labels is generally smoother, yet TITRL still shows clear advantages across all datasets, which confirms the effectiveness of the label semantic learning in our method. Thus, TITRL maintains stable and robust performance under both feature absence and label incompleteness, highlighting its applicability to challenging real-world scenarios.

To provide a more comprehensive multi-metric visualization of performance, we further present radar plots under various conditions of data sparsity. Fig. 8 illustrates the results with a fixed LMR of 50% while PER increases from 50% to 90%. It is evident that TITRL consistently occupies the outermost boundary across all datasets and PER levels. This visually confirms its superior performance across all six evaluation metrics. Even as the view missingness becomes more severe (e.g., PER=90%), TITRL maintains a significant performance margin over the competing methods, highlighting its robustness against feature absence.

Furthermore, Fig. 9 and Fig. 10 depict scenarios with even higher label missingness (LMR=70% and LMR=90%, respectively). In these highly challenging settings, TITRL’s superiority remains unshaken. As data incompleteness intensifies, the performance of many baseline methods degrades substantially, whereas TITRL sustains its strong performance profile. This demonstrates the exceptional resilience of our method in handling extreme levels of label sparsity. Collectively, these radar plots offer strong visual evidence that TITRL not only achieves the highest performance but also maintains remarkable robustness, solidifying its effectiveness for real-world iMvMLC tasks.

##### Parameter Sensitivity.

To investigate the sensitivity of the model’s hyperparameters, we perform a grid search for parameter selection on  $\lambda_1$  and  $\lambda_2$ , and report the validation AP on Corel5k, ESPGame, and IAPRTC12. The results are visualized as heatmaps in Fig. 6. The values of both parameters are adjusted within the range  $\{0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1\}$ . It can be observed that TITRL is more responsive to  $\lambda_1$ , as an improper selection of  $\lambda_1$  can lead to a significant performance degradation. The optimal configuration tends to appear when  $\lambda_1$  lies within the range (0.01, 0.05) while  $\lambda_2$  approaches 1, a trend that holds consistently across both small-scale datasets (e.g., Corel 5k, IAPRTC12) and large-scale datasets (e.g., Mirflickr, OBJECT). These observations suggest that TITRL requires careful tuning of  $\lambda_1$  to promote the semantic compactness of the extracted representations, whereas  $\lambda_2$  can be set to a relatively high value to stabilize performance. Overall, TITRL maintains stable accuracy across a wide parameter range, with clearly identifiable regions of optimality.

##### Convergence Behavior.

We further examine the optimization behavior of TITRL by tracking training and validation dynamics over epochs. Fig. 7 presents the evolution curves of training loss and validation AP on Corel5k,

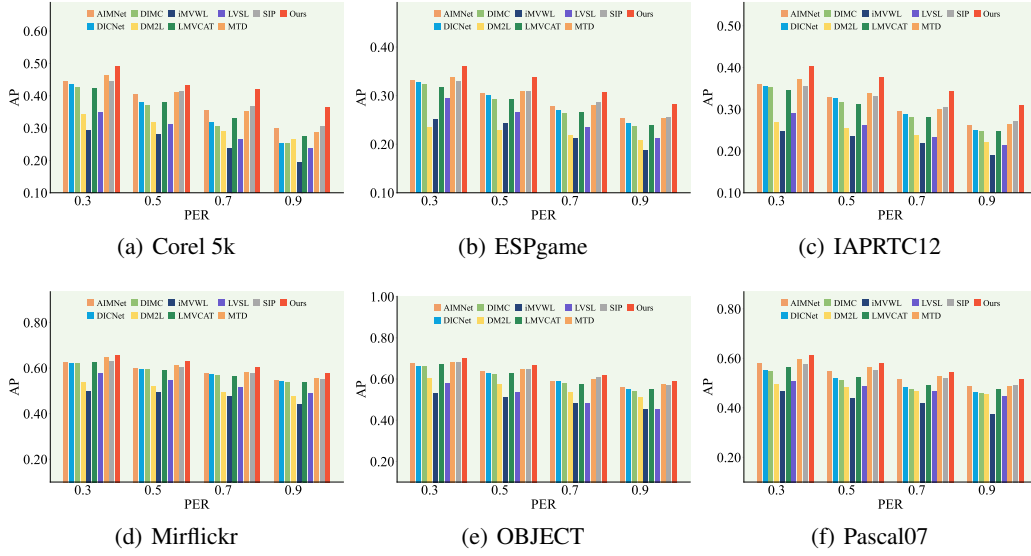


Figure 4: Results on six datasets with PER changing from 30% to 90%, and LMR fixed at 50%

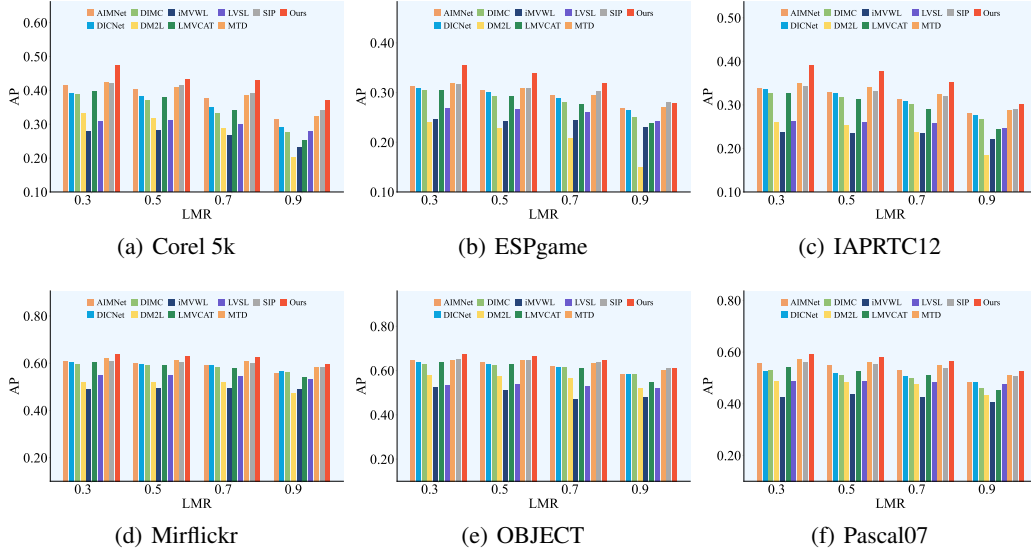


Figure 5: Results on six datasets with LMR changing from 30% to 90%, and PER fixed at 50%

ESPGAME, and IAPRTC12. From the result, we can find that the training loss decreases smoothly, while validation AP rises quickly in the early epochs and stabilizes thereafter, typically converging within 25-40 epochs. This rapid and stable convergence demonstrates the efficiency of our optimization strategy and the well-conditioned nature of the proposed objective. Moreover, the close alignment between training and validation curves indicates that overfitting is effectively controlled. We also observe consistent convergence patterns across datasets and repeated runs, which confirms the numerical stability of TITRL under diverse conditions. These properties collectively highlight that TITRL not only achieves strong predictive accuracy under missing data but also provides reliable training dynamics, ensuring efficiency and reproducibility in practice.

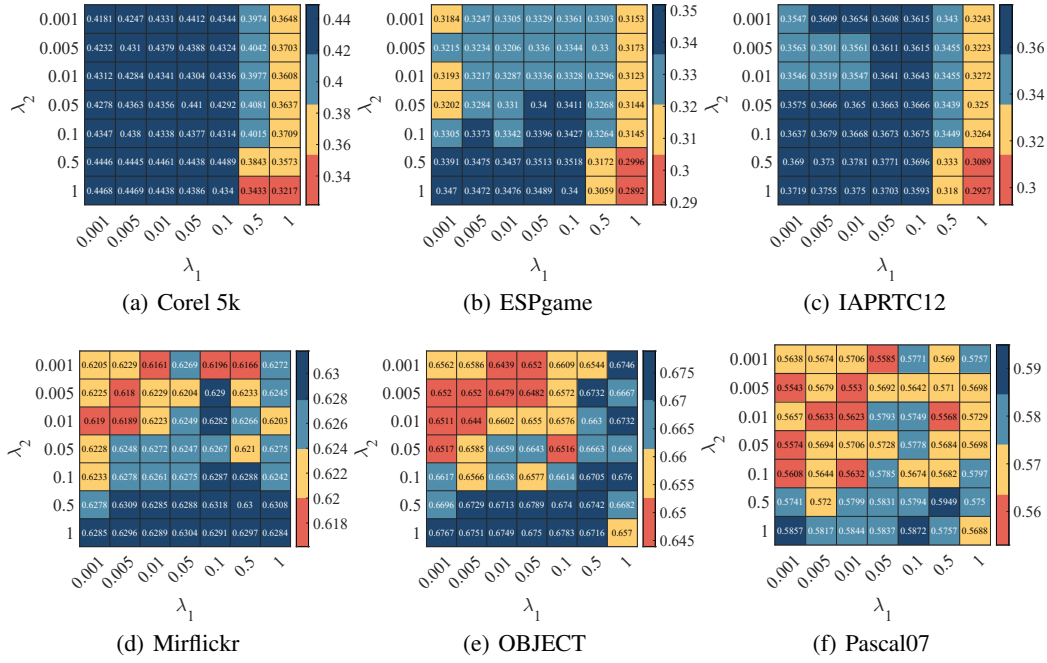


Figure 6: Parameter sensitivity analysis.

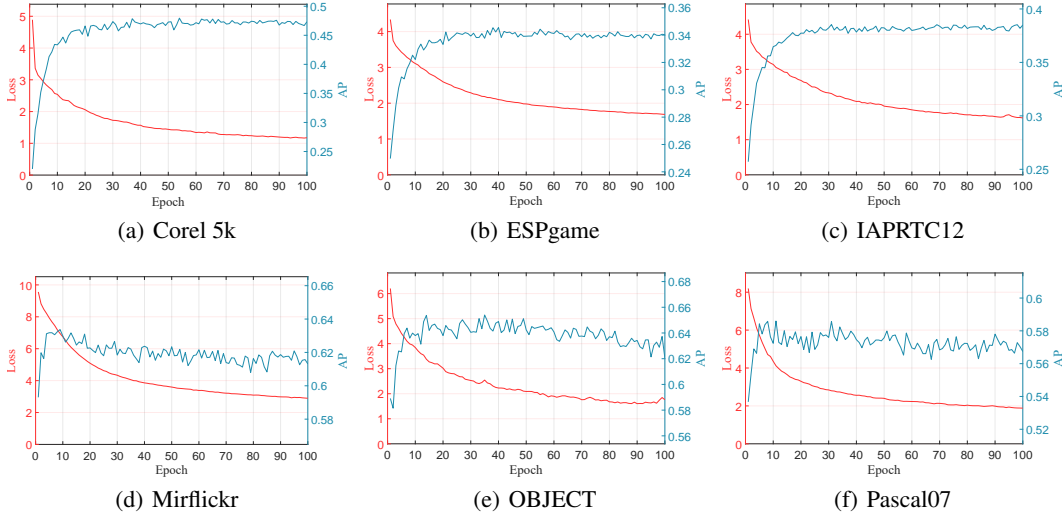


Figure 7: Convergence analysis.

## USE OF LARGE LANGUAGE MODELS (LLMs)

In the course of this research, Large Language Models (LLMs) were employed as an assistive tool for light English editing, such as grammar corrections, wording improvements, and enhancing clarity. The LLM played no role in the ideation, experimental design, data processing, analysis, or methodological development of the research. All technical content, including research ideas, methodology, data analysis, and results, were independently developed and verified by the authors, who take full responsibility for the manuscript.

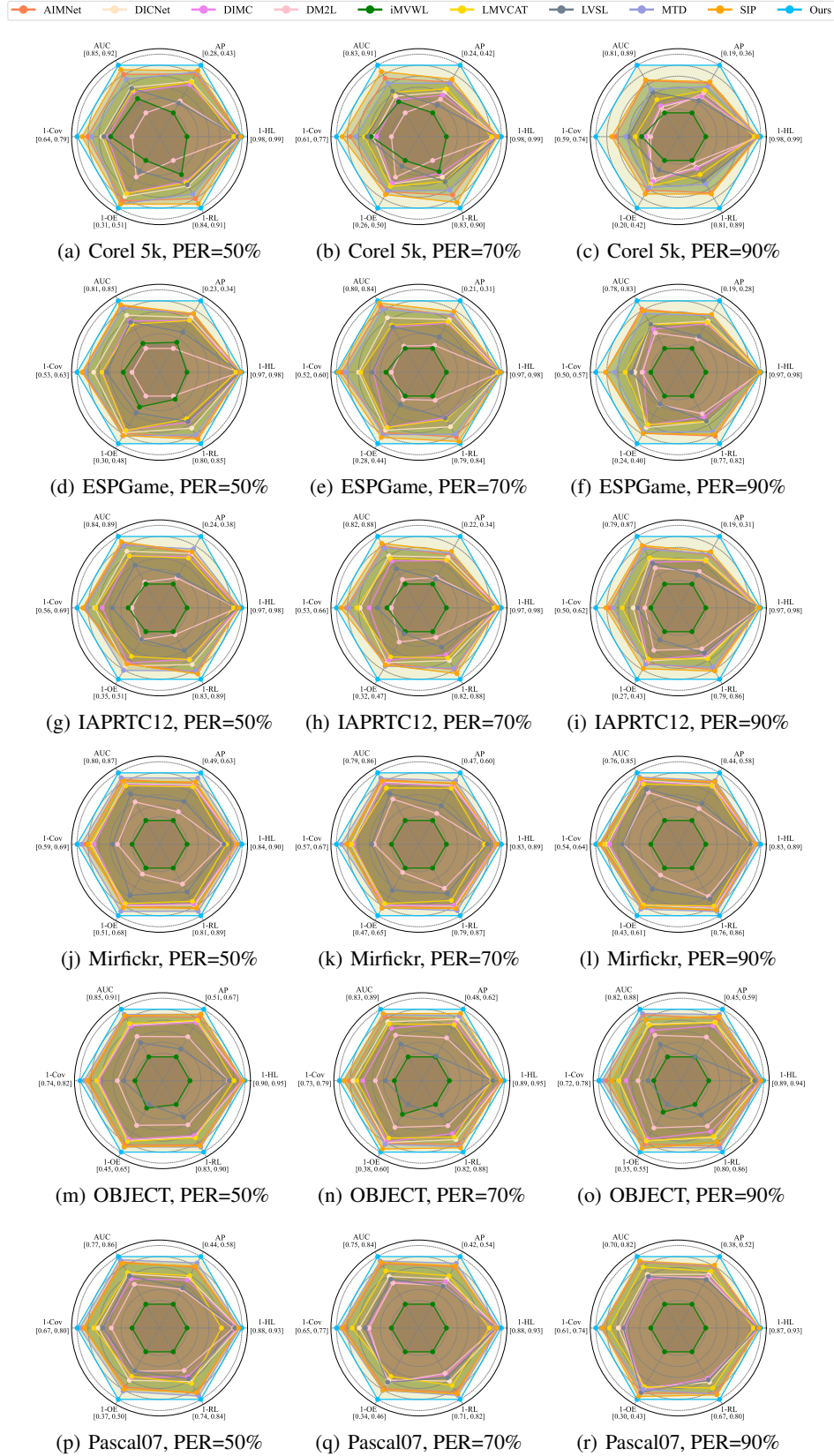


Figure 8: Experimental results of ten methods on six datasets with PER varying from 50% to 90%, while LMR = 50%.

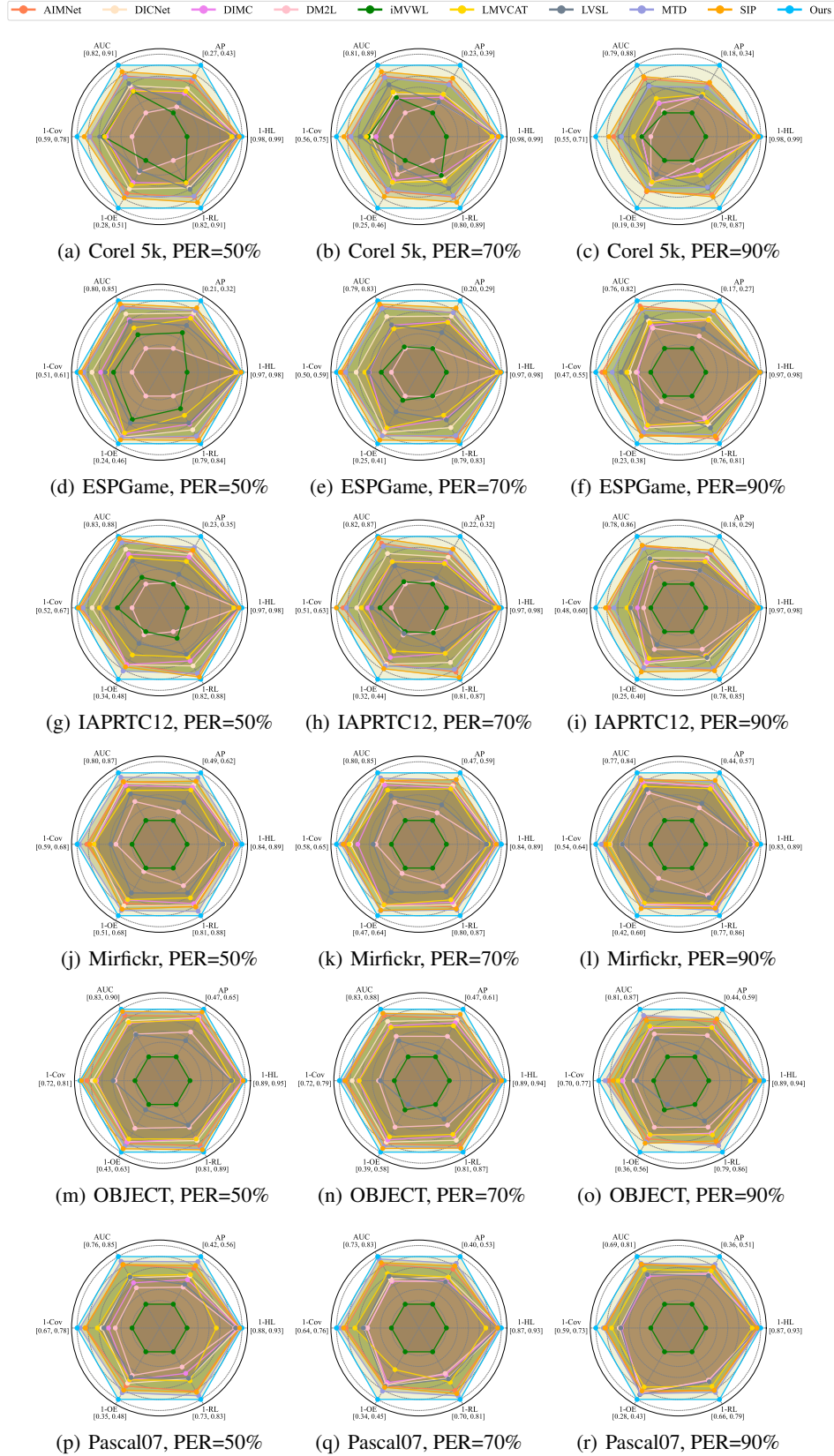


Figure 9: Experimental results of ten methods on six datasets with PER varying from 50% to 90%, while LMR=70%.



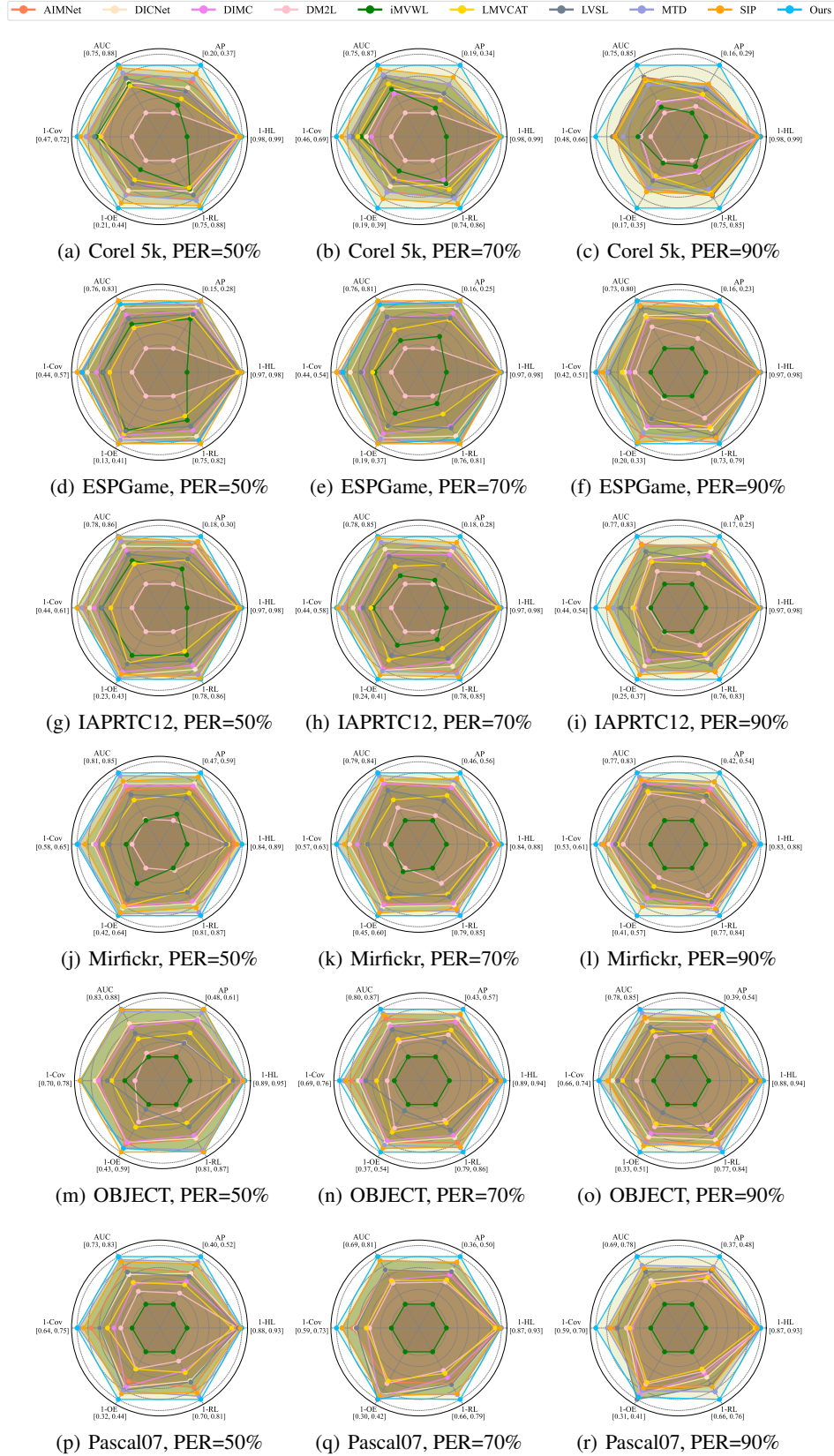


Figure 10: Experimental results of ten methods on six datasets with PER varying from 50% to 90%, while LMR = 90%.