

Theory-Inspired Task-Relevant Representation Learning for Incomplete Multi-View Multi-Label Learning

Anonymous CVPR submission

Abstract

001 *Multi-view multi-label learning is commonly hindered by*
002 *dual data incompleteness, arising from constraints in fea-*
003 *ture collection and prohibitive annotation costs. To ad-*
004 *dress the intricate yet highly practical challenges and en-*
005 *hance the reliability of representation extraction, hetero-*
006 *geneous feature fusion, and label semantic learning, we*
007 *propose a Theory-Inspired Task-Relevant Representation*
008 *Learning method named TITRL. From an information-*
009 *theoretic standpoint, we identify the sources of view-specific*
010 *information that interfere with shared representations. By*
011 *introducing dual-layer constraints on feature exclusivity*
012 *and label integration, TITRL constructs a general frame-*
013 *work for task-relevant information extraction. Besides,*
014 *through variational derivation, we demonstrate the exis-*
015 *tence of tractable bounds for the mutual information model*
016 *that guides the optimization direction. Regarding label se-*
017 *mantic learning, we establish flexible relationships between*
018 *label prototypes by promoting the expression of sample-*
019 *level label correlations. During the multi-view integration*
020 *process , TITRL simultaneously incorporates early fusion*
021 *through distribution information aggregation and late fu-*
022 *sion weighted by prediction confidence, which improves the*
023 *semantic stability while enabling dynamic view quality as-*
024 *essment. Finally, extensive experimental results validate*
025 *the effectiveness of TITRL against state-of-the-art methods.*

026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077

ously tagged with labels such as “forest”, “urban area”, and “water body”[25]. 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[23, 28]. 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 [19, 31].

Existing MvMLC techniques seek to leverage heterogeneous features and predict multiple labels within a unified framework. Representative methods include lrMMC [20] applying low-rank matrix factorization, and EF²FS [11] 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 [9]. 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 [26]. 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 [33] and supervised classification tasks [2]. DICNet [15] and LMVCAT [16] 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 [18] 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 [30] and AIMNet [17], [17, 30] 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 [14], 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 [24]. For example, in a movie recommendation system [13], 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 prediction confidence 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

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 [7] have demonstrated that integrating the common informa-

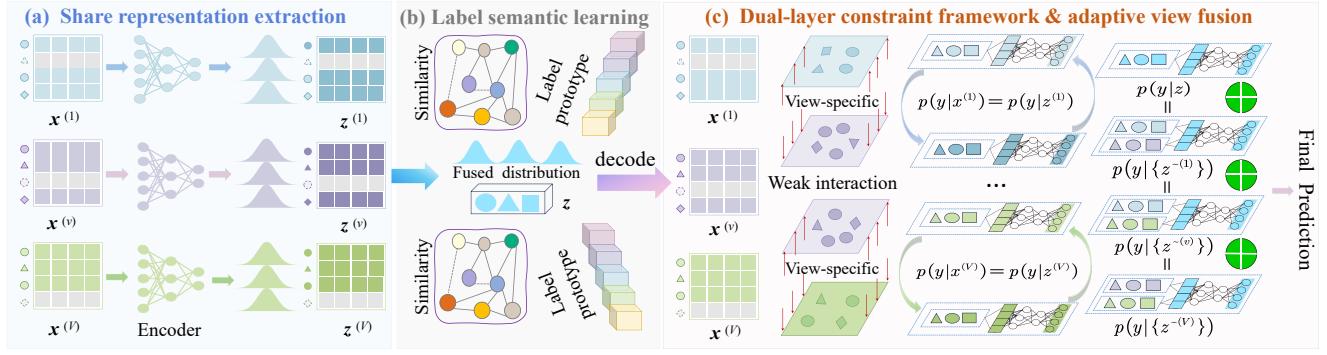


Figure 1. The main framework of our proposed TITRL. Different shapes signify different samples.

tion and reducing the redundant information introduced by view-specific factors is sufficient to accomplish all prediction tasks. For example, in facial recognition [21], 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 and background clutter, which can disrupt model performance. By integrating shared features and removing noisy information, the model can enhance recognition accuracy. 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 aggregating information from 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 $I(\{\mathbf{x}^{\sim(v)}\}; \mathbf{z}^{(v)} | \mathbf{x}^{(v)})$, where $\{\{\mathbf{x}^{\sim(v)}\}, \mathbf{x}^{(v)}\} = \{\mathbf{x}\}$. Next, due to the scalability of information transfer, we can derive the following upper bound to guide information separation:

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

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:

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

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 rep-

245 representations is obtained:

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

246 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:

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

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

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

255 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:

$$I(\{\mathbf{x}^{\sim(v)}\}; \mathbf{z} | \mathbf{x}^{(v)}) \\ = \iint p(\{\mathbf{x}\}, \mathbf{z}) \log \frac{p(\mathbf{z} | \{\mathbf{x}\})}{g^v(\mathbf{z} | \mathbf{x}^{(v)})} d\{\mathbf{x}\} d\mathbf{z} \\ - \int p(\mathbf{x}^{(v)}) D_{KL} \left(p(\mathbf{z} | \mathbf{x}^{(v)}) \| g^v(\mathbf{z} | \mathbf{x}^{(v)}) \right) d\mathbf{x}^{(v)}. \quad (7)$$

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

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

266 In conclusion, the trainable loss subject to the feature-level constraint is given by

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

270 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

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

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

$$I(\mathbf{y}; \mathbf{z}^{(v)} | \{\mathbf{z}^{\sim(v)}\}) = H(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) - H(\mathbf{y} | \{\mathbf{z}\}) \\ = - \int p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \log p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) d\mathbf{y} \\ + \int p(\mathbf{y} | \{\mathbf{z}\}) \log p(\mathbf{y} | \{\mathbf{z}\}) d\mathbf{y}. \quad (11)$$

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

$$I(\mathbf{y}; \mathbf{z}^{(v)} | \{\mathbf{z}^{\sim(v)}\}) \leq D_{KL} \left(p(\mathbf{y} | \{\mathbf{z}\}) \| p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}) \right) \\ + H \left(p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\}), p(\mathbf{y} | \{\mathbf{z}\}) \right). \quad (12)$$

282 Owing to the congruent optimization objective of aligning $p(\mathbf{y} | \{\mathbf{z}\})$ and $p(\mathbf{y} | \{\mathbf{z}^{\sim(v)}\})$, we directly 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 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

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

2.3. Adaptive View Fusion and Label Representation Learning

291 The integration of view representations constitutes a critical challenge in multi-view learning. Given that variational inference optimizes the distribution of each view, we 293 leverage the progressively refined distribution information 295 to facilitate view fusion. Within the network architecture, 296 each view is processed through two encoders to estimate 297 the latent distribution of its shared representations. Specif- 299 ically, the distribution is modeled as $p(\mathbf{z}^{(v)} | \mathbf{x}^{(v)}) := 300 \mathcal{N}(f_\mu^v(\mathbf{x}^{(v)}), f_{\sigma^2}^v(\mathbf{x}^{(v)}) \mathbf{I})$, where f_μ^v and $f_{\sigma^2}^v$ 301 are the mean and variance encoders. To ensure that the shared feature 302 incorporates information from all views and benefits from 303

304 the greater stability of representations with lower variance,
 305 we adopt the product-of-experts (PoE) framework (Hinton,
 306 2002) with the one-vote property to perform weighted fu-
 307 sion of the distribution parameters across views:

$$308 \quad \begin{cases} \mu_s = \frac{\sum_{v \in \mathcal{V}} f_\mu^v(\mathbf{x}^{(v)}) \frac{1}{f_{\sigma^2}^v(\mathbf{x}^{(v)})}}{\sum_{v \in \mathcal{V}} \frac{1}{f_{\sigma^2}^v(\mathbf{x}^{(v)})} + 1} \\ \sigma_s^2 = \frac{1}{\sum_{v \in \mathcal{V}} \frac{1}{f_{\sigma^2}^v(\mathbf{x}^{(v)})} + 1}. \end{cases} \quad (14)$$

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

$$311 \quad \mathbf{z} = \frac{1}{S} \sum_{i=1}^S (\mu_s + \sigma_s \odot \delta^i), \quad (15)$$

312 where $\delta^i \in \mathbb{R}^d$ denotes the i -th sampling from the standard
 313 Gaussian distribution and \odot indicates element-wise mul-
 314 tiplication. The representations $\{\mathbf{z}^{(v)}\}_{v=1}^V$ extracted from
 315 each view are also sampled from their respective distribu-
 316 tions following Eq. (15). During view fusion, it is essential
 317 to not only account for the aggregation of representation
 318 information but also to incorporate the impact of label in-
 319 formation. Since multiple labels are typically encoded as
 320 one-hot vectors, which lacks the flexibility to capture la-
 321 bel semantics, particularly in scenarios with missing labels.
 322 To address this, we adopt a data-driven approach to intro-
 323 duce label prototypes, ensuring that the semantic informa-
 324 tion carried by these prototypes is closely aligned with the
 325 ground truth labels. In order to explicitly model label ex-
 326 pressions, we employ stochastic encoders to fit the underly-
 327 ing distribution $\mathcal{N}(\mu_i, \sigma_i^2 \mathbf{I})$ for each label prototype, where
 328 μ_i and σ_i^2 are the d -dimensional mean and variance outputs,
 329 respectively, produced by the encoders $h_\mu(\mathbf{b}_i)$ and $h_{\sigma^2}(\mathbf{b}_i)$.
 330 $\mathbf{b}_i \in \mathbb{R}^C$ serves as a learnable embedding corresponding to
 331 the i -th class, which is initialized as a one-hot vector with
 332 the i -th entry is 1. After obtaining $\{\mathbf{l}_i\}_{i=1}^C$ through stochas-
 333 tic sampling, it is necessary to capture the intricate corre-
 334 lations between these label representations, which forms a
 335 crucial determinant in enhancing the performance of multi-
 336 label classification. Considering that the manifestation of
 337 label correlations differs across samples, we adopt a nu-
 338 nanced approach that centers on instance-level relevance to
 339 strengthen the similarity between the cross-view represen-
 340 tation of each sample and the label attributes it possesses.
 341 Specifically, we sample the shared feature \mathbf{z} according to
 342 Eq. (15), with its associated known label prototype being
 343 $\{\mathbf{l}_i | i \in \mathcal{U}\}$. By using the cosine similarity as the criterion,
 344 the alignment loss designed to capture label correlations is

$$345 \quad \mathcal{L}_a = \frac{1}{|\mathcal{U}|} \sum_{i \in \mathcal{U}} \frac{\langle \mathbf{z} \cdot \mathbf{l}_i \rangle}{\|\mathbf{z}\| \|\mathbf{l}_i\|}. \quad (16)$$

By optimizing loss \mathcal{L}_a , we refine the mapping semantic be-
 346 between features and labels, while simultaneously highlight-
 347 ing the associations among label prototypes, which are tai-
 348 lored for application to each individual sample. Subse-
 349 quently, during label prediction, it is important to synthe-
 350 size the generalized information from multiple views with
 351 the semantic representations of individual categories. When
 352 these information exhibit coherence, it becomes feasible to
 353 infer that the sample contains the relevant labels. To this
 354 end, we utilize a neural network to adaptively gauge the
 355 degree of similarity between view representations and cate-
 356 gory 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 σ_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 $p(\mathbf{y}|\{\mathbf{z}^{(v)}\})$ 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 $p(\mathbf{y}|\{\mathbf{z}^{(v)}\})$. 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:

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

where all weighting coefficients $\mathcal{L}_{con}^{(v)} (0 \leq v \leq V)$ are
 384 normalized in advance. To enhance the classification dis-
 385 criminating power and reinforce the interaction term $I(\mathbf{y}; \mathbf{z}^{(v)})$
 386 in model (4), we employ the following cross-entropy loss:
 387

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

The classification loss in our method is the aggrega-
 389 tion of four distinct components, with one arising from
 390 the final prediction and the remaining three emanating
 391

392 from pseudo-predictions $p(\mathbf{y}_i|\{\mathbf{z}^{\sim(v)}\})$, $p(\mathbf{y}_i|\mathbf{x}^{(v)})$, and
 393 $p(\mathbf{y}_i|\mathbf{z}^{(v)})$, which collectively constitutes the overall loss
 394 \mathcal{L}_{BCE}^t . Thus, the total training loss of TITRR is as below:

$$395 \quad \mathcal{L} = \mathcal{L}_{BCE}^t + \mathcal{L}_a + \lambda_1 \mathcal{L}_f + \lambda_2 \mathcal{L}_c, \quad (20)$$

396 where λ_1 and λ_2 govern the trade-off between the empirical
 397 values and impacts of different losses.

398 3. Experiments

399 3.1. Datasets and Metrics

400 In our experiments, we employ six widely used multi-
 401 view multi-label datasets to evaluate the effectiveness of
 402 our method, i.e., Corel 5k [3], ESPGame [1], IAPRTC12
 403 [8], Mirflickr [12], Pascal07 [4], and OBJECT [10]. Fol-
 404 lowing the evaluation protocols in [15, 30], we adopt the
 405 following six metrics to form a comprehensive assessment
 406 framework, i.e., Hamming Loss (HL), Ranking Loss (RL),
 407 OneError (OE), Coverage (Cov), Average Precision (AP),
 408 and Area Under Curve (AUC). For clarity in comparison,
 409 we report 1-HL, 1-OE, 1-Cov, and 1-RL, where higher val-
 410 ues consistently indicate better performance.

411 3.2. Comparison Methods

412 To assess the performance of our method, we compare it
 413 with nine state-of-the-art approaches, i.e., AIMNet [17],
 414 DICNet [15], DIMC [30], iMVWL [27], LMVCAT [16],
 415 MTD [14], SIP [18], LVSL [32], and DM2L [22]. The
 416 first seven methods are capable of simultaneously handling
 417 missing views and labels. Since LVSL cannot directly pro-
 418 cess incomplete data, we impute missing views using the
 419 mean of available instances and fill absent labels with ze-
 420 ros. DM2L is a kernel-based nonlinear method for incom-
 421 plete multi-label learning. Thus, we concatenate all recov-
 422 ered views into a single representation to apply DM2L. All
 423 hyperparameters of the compared methods are set according
 424 to the recommended configurations in their original imple-
 425 ments, ensuring a fair and reproducible comparison.

426 3.3. Implementation Details

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

437 3.4. Experimental Results and Analysis

438 To rigorously evaluate the effectiveness of TITRL in han-
 439 dling absent views and incomplete labels, we conduct ex-

440 tensive comparative experiments against nine representa-
 441 tive algorithms across six benchmark datasets under vary-
 442 ing levels of data sparsity. Specifically, the proportions
 443 of missing views (PER) and labels (LMR) are set to
 444 $\{30\%, 50\%, 70\%, 90\%\}$. The results in terms of the mean
 445 and standard deviation at PER=50% and LMR=50% are
 446 summarized in Table 1, along with the average ranking
 447 across six evaluation metrics to provide an aggregated per-
 448 formance assessment. In addition, Fig. 2 visualizes how AP
 449 evolves as the missing proportion increases, while Fig. 3
 450 presents radar plots that jointly capture multi-metric per-
 451 formance distribution at PER=90% and LMR=90%. These re-
 452 sults collectively provide a holistic evaluation of predictive
 453 accuracy and model robustness.

454 From the comparison results, several important observa-
 455 tions can be drawn: (i) TITRL consistently secures the best
 456 results across almost all datasets and metrics. For instance,
 457 on Corel5k, TITRL achieves an AP score of 0.432, outper-
 458 forming SIP (0.414) and MTD (0.410), with similar margins
 459 observed on other datasets. Besides, TITRL maintains its
 460 superiority on large-scale datasets like ESPGame and Mir-
 461 flickr, which underscores its scalability and resilience. (ii)
 462 Drawn from Fig. 2, we can find that while competing meth-
 463 ods suffer steep performance degradation when missing
 464 ratios achieve a high level, TITRL continues to exhibit con-
 465 siderable performance. For example, when PER=70%, the
 466 performance of all comparison methods stays below 0.36,
 467 while TITRL surpasses 0.4 by a certain margin. Although
 468 our method still outperforms others under conditions of se-
 469 vere label missingness, the performance gap is prominently
 470 reflected in the presence of feature absent. This further
 471 highlights the critical importance of multi-view represen-
 472 tation learning, a role that our method is well equipped to ful-
 473 fill. As depicted in the radar chart of Fig. 3, it is evident that
 474 our method consistently occupies the outermost boundary,
 475 which indicates that TITRL stands out even under highly
 476 challenging conditions across various evaluation perspec-
 477 tives. Consequently, our method demonstrates strong ro-
 478 bustness in addressing the problem of data incompleteness
 479 and shows significant potential for broader adoption. (iii)
 480 As evidenced by Table 1, our method consistently main-
 481 tains the top position, while the rankings of alternative ap-
 482 proaches remain volatile, which shows the exceptional per-
 483 formance stability of our method. Against top competing
 484 methods SIP and MTD, our approach also demonstrates su-
 485 periority, which underscores the pivotal role of introducing
 486 label integration strategies in the feature extraction process
 487 and fine-grained characterization of label semantics. Com-
 488 compared with the deep learning-based methods AIMNet, DIC-
 489 Net and DIMC, the substantial advantage of TITRL further
 490 reveals the importance of jointly considering the view prop-
 491 erty 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 _{0.000}	0.987 _{0.000}	0.987 _{0.000}	0.987 _{0.000}	0.978 _{0.000}	0.986 _{0.000}	0.987 _{0.000}	0.988 _{0.000}	0.988 _{0.000}	0.988 _{0.000}
	1-OE	0.478 _{0.011}	0.460 _{0.012}	0.446 _{0.009}	0.378 _{0.014}	0.308 _{0.017}	0.448 _{0.011}	0.353 _{0.017}	0.492 _{0.011}	0.492 _{0.014}	0.509 _{0.014}
	1-Cov	0.766 _{0.004}	0.726 _{0.007}	0.709 _{0.008}	0.640 _{0.007}	0.701 _{0.003}	0.720 _{0.006}	0.720 _{0.005}	0.754 _{0.005}	0.781 _{0.004}	0.795 _{0.006}
	1-RL	0.900 _{0.002}	0.881 _{0.004}	0.874 _{0.004}	0.843 _{0.004}	0.864 _{0.002}	0.876 _{0.004}	0.879 _{0.002}	0.893 _{0.004}	0.908 _{0.003}	0.914 _{0.003}
	AP	0.404 _{0.005}	0.381 _{0.006}	0.370 _{0.005}	0.318 _{0.005}	0.281 _{0.005}	0.379 _{0.006}	0.311 _{0.005}	0.410 _{0.007}	0.414 _{0.006}	0.432 _{0.007}
	AUC	0.903 _{0.002}	0.883 _{0.004}	0.877 _{0.004}	0.846 _{0.004}	0.867 _{0.002}	0.879 _{0.003}	0.882 _{0.002}	0.896 _{0.003}	0.910 _{0.002}	0.916 _{0.002}
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 _{0.000}	0.983 _{0.000}	0.983 _{0.000}	0.983 _{0.000}	0.972 _{0.000}	0.982 _{0.000}	0.983 _{0.000}	0.983 _{0.000}	0.983 _{0.000}	0.983 _{0.000}
	1-OE	0.442 _{0.006}	0.440 _{0.009}	0.431 _{0.009}	0.302 _{0.008}	0.343 _{0.010}	0.432 _{0.006}	0.365 _{0.006}	0.452 _{0.007}	0.450 _{0.006}	0.481 _{0.006}
	1-Cov	0.621 _{0.003}	0.601 _{0.003}	0.586 _{0.004}	0.532 _{0.003}	0.548 _{0.004}	0.587 _{0.003}	0.578 _{0.002}	0.617 _{0.004}	0.622 _{0.004}	0.631 _{0.008}
	1-RL	0.845 _{0.002}	0.836 _{0.002}	0.830 _{0.002}	0.804 _{0.002}	0.807 _{0.002}	0.827 _{0.002}	0.829 _{0.001}	0.843 _{0.002}	0.847 _{0.002}	0.852 _{0.004}
	AP	0.306 _{0.003}	0.300 _{0.003}	0.294 _{0.003}	0.229 _{0.003}	0.243 _{0.004}	0.293 _{0.003}	0.266 _{0.003}	0.309 _{0.003}	0.309 _{0.004}	0.339 _{0.003}
	AUC	0.850 _{0.001}	0.841 _{0.002}	0.835 _{0.002}	0.808 _{0.001}	0.813 _{0.002}	0.832 _{0.001}	0.834 _{0.001}	0.847 _{0.002}	0.851 _{0.002}	0.855 _{0.003}
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 _{0.000}	0.981 _{0.000}	0.981 _{0.000}	0.980 _{0.000}	0.969 _{0.000}	0.980 _{0.000}	0.981 _{0.000}	0.981 _{0.000}	0.981 _{0.000}	0.982 _{0.000}
	1-OE	0.457 _{0.008}	0.464 _{0.008}	0.454 _{0.006}	0.378 _{0.008}	0.351 _{0.008}	0.433 _{0.009}	0.377 _{0.007}	0.479 _{0.007}	0.459 _{0.005}	0.508 _{0.008}
	1-Cov	0.675 _{0.004}	0.649 _{0.005}	0.630 _{0.005}	0.556 _{0.005}	0.565 _{0.004}	0.646 _{0.004}	0.605 _{0.004}	0.670 _{0.004}	0.678 _{0.003}	0.693 _{0.006}
	1-RL	0.884 _{0.001}	0.874 _{0.002}	0.868 _{0.002}	0.837 _{0.002}	0.833 _{0.002}	0.868 _{0.002}	0.857 _{0.002}	0.882 _{0.002}	0.886 _{0.001}	0.893 _{0.002}
	AP	0.329 _{0.003}	0.326 _{0.003}	0.318 _{0.002}	0.254 _{0.002}	0.236 _{0.002}	0.313 _{0.004}	0.262 _{0.002}	0.340 _{0.002}	0.331 _{0.002}	0.377 _{0.004}
	AUC	0.885 _{0.001}	0.876 _{0.002}	0.870 _{0.001}	0.838 _{0.001}	0.835 _{0.001}	0.870 _{0.002}	0.859 _{0.001}	0.883 _{0.002}	0.887 _{0.001}	0.894 _{0.002}
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 _{0.001}	0.890 _{0.001}	0.890 _{0.001}	0.876 _{0.001}	0.840 _{0.004}	0.880 _{0.004}	0.877 _{0.001}	0.893 _{0.000}	0.890 _{0.001}	0.896 _{0.001}
	1-OE	0.646 _{0.009}	0.647 _{0.010}	0.646 _{0.008}	0.533 _{0.008}	0.511 _{0.016}	0.639 _{0.009}	0.609 _{0.007}	0.667 _{0.006}	0.654 _{0.007}	0.683 _{0.006}
	1-Cov	0.673 _{0.003}	0.662 _{0.004}	0.657 _{0.003}	0.615 _{0.002}	0.588 _{0.013}	0.665 _{0.002}	0.624 _{0.002}	0.681 _{0.002}	0.669 _{0.006}	0.688 _{0.003}
	1-RL	0.874 _{0.002}	0.869 _{0.003}	0.867 _{0.003}	0.835 _{0.001}	0.809 _{0.014}	0.862 _{0.003}	0.847 _{0.001}	0.878 _{0.001}	0.873 _{0.002}	0.886 _{0.002}
	AP	0.599 _{0.003}	0.595 _{0.007}	0.592 _{0.006}	0.519 _{0.003}	0.495 _{0.017}	0.589 _{0.004}	0.548 _{0.003}	0.614 _{0.004}	0.603 _{0.005}	0.629 _{0.004}
	AUC	0.861 _{0.001}	0.855 _{0.002}	0.854 _{0.002}	0.828 _{0.001}	0.801 _{0.017}	0.852 _{0.003}	0.839 _{0.001}	0.864 _{0.001}	0.859 _{0.002}	0.871 _{0.002}
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 _{0.001}	0.948 _{0.001}	0.947 _{0.001}	0.935 _{0.000}	0.899 _{0.002}	0.940 _{0.002}	0.935 _{0.001}	0.949 _{0.001}	0.948 _{0.001}	0.950 _{0.001}
	1-OE	0.619 _{0.015}	0.601 _{0.011}	0.594 _{0.012}	0.537 _{0.011}	0.465 _{0.018}	0.604 _{0.016}	0.450 _{0.008}	0.627 _{0.011}	0.626 _{0.009}	0.648 _{0.008}
	1-Cov	0.807 _{0.006}	0.794 _{0.006}	0.793 _{0.006}	0.768 _{0.005}	0.744 _{0.008}	0.796 _{0.008}	0.759 _{0.006}	0.813 _{0.006}	0.809 _{0.006}	0.818 _{0.006}
	1-RL	0.888 _{0.005}	0.876 _{0.004}	0.875 _{0.004}	0.860 _{0.004}	0.833 _{0.006}	0.878 _{0.006}	0.850 _{0.004}	0.890 _{0.005}	0.889 _{0.004}	0.897 _{0.003}
	AP	0.639 _{0.010}	0.627 _{0.009}	0.623 _{0.010}	0.577 _{0.008}	0.512 _{0.014}	0.630 _{0.012}	0.537 _{0.008}	0.649 _{0.009}	0.649 _{0.008}	0.665 _{0.006}
	AUC	0.897 _{0.004}	0.886 _{0.004}	0.885 _{0.004}	0.872 _{0.004}	0.846 _{0.006}	0.888 _{0.006}	0.864 _{0.004}	0.900 _{0.005}	0.898 _{0.004}	0.906 _{0.003}
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 _{0.001}	0.931 _{0.000}	0.931 _{0.001}	0.927 _{0.001}	0.882 _{0.004}	0.915 _{0.005}	0.928 _{0.001}	0.933 _{0.001}	0.932 _{0.001}	0.935 _{0.001}
	1-OE	0.462 _{0.010}	0.443 _{0.007}	0.435 _{0.010}	0.419 _{0.006}	0.366 _{0.039}	0.433 _{0.016}	0.418 _{0.008}	0.474 _{0.008}	0.468 _{0.008}	0.496 _{0.009}
	1-Cov	0.781 _{0.007}	0.749 _{0.003}	0.738 _{0.010}	0.720 _{0.004}	0.674 _{0.011}	0.759 _{0.006}	0.738 _{0.003}	0.790 _{0.006}	0.778 _{0.004}	0.795 _{0.005}
	1-RL	0.830 _{0.006}	0.804 _{0.002}	0.792 _{0.008}	0.778 _{0.003}	0.736 _{0.011}	0.808 _{0.006}	0.797 _{0.002}	0.836 _{0.006}	0.828 _{0.004}	0.844 _{0.004}
	AP	0.549 _{0.007}	0.517 _{0.004}	0.510 _{0.008}	0.482 _{0.005}	0.438 _{0.022}	0.524 _{0.009}	0.486 _{0.005}	0.562 _{0.005}	0.552 _{0.006}	0.581 _{0.007}
	AUC	0.851 _{0.005}	0.827 _{0.002}	0.817 _{0.008}	0.806 _{0.003}	0.767 _{0.011}	0.830 _{0.006}	0.823 _{0.002}	0.855 _{0.006}	0.848 _{0.005}	0.861 _{0.003}
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					
			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	
✓	✗	✓	0.561	0.845	0.834	0.476	0.652	0.896	0.889	0.632	
✓	✓	✗	0.584	0.861	0.847	0.495	0.668	0.905	0.897	0.646	
✓	✓	✓	0.588	0.867	0.851	0.501	0.673	0.910	0.901	0.651	

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 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 instrumen-

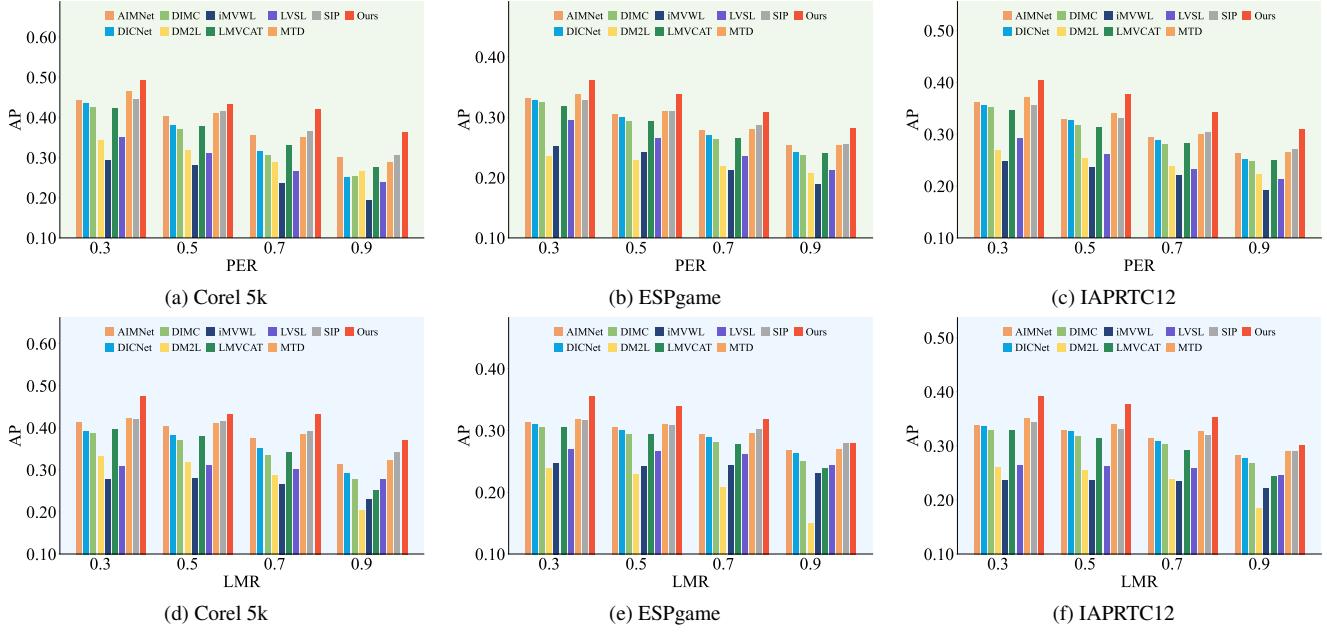


Figure 2. Experimental results on three datasets with one of PER and LMR fixed at 50% and the other varying from 30% to 90%

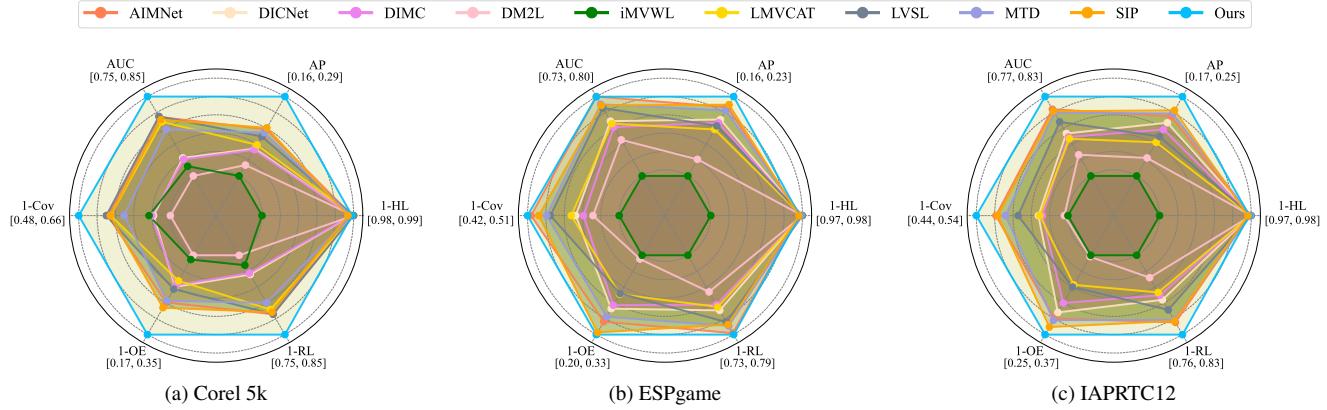


Figure 3. Experimental results of ten methods on three datasets with PER= 90% and LMR= 90%.

507
508
509
510
511
512
tal 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.

513 4. Conclusion

514
515
516
517
518
519
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 common information, improve the reliability of view
fusion, and accurately capture label correlation semantics.
Specifically, TITRL introduces a dual-layer mutual infor-

520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
mation constraint framework, which enables the disentangle-
ment of view-specific noise. By deriving tractable vari-
ational 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 property with a confidence-driven
late fusion mechanism, thereby enhancing the stability of
representation expression. Moreover, we explicitly model
sample-level label correlations by aligning shared represen-
tations with learnable label prototypes, allowing for flexible
use of label dependencies. Extensive experiments demon-
strate TITRL’s superiority, particularly in high-deficiency
conditions where most baselines fail. In the future, we plan
to utilize the prior knowledge embedded in large language
models to jointly learn the semantics at the feature and label
levels.

536 References

- 537 [1] Luis Von Ahn and Laura Dabbish. Labeling images with a
538 computer game. In *SIGCHI Conference on Human Factors
539 in Computing Systems*, page 319–326, 2004. 6
- 540 [2] Wei Chen, Jiage Chen, Yuewu Wan, Xining Liu, Mengya
541 Cai, Jingguo Xu, Hongbo Cui, and Mengdie Duan. Land
542 cover classification based on multimodal remote sensing fu-
543 sion. *ISPRS Annals of the Photogrammetry, Remote Sensing
544 and Spatial Information Sciences*, 10:35–40, 2024. 2
- 545 [3] Pinar Duygulu, Kobus Barnard, Joao FG de Freitas, and
546 David A Forsyth. Object recognition as machine transla-
547 tion: Learning a lexicon for a fixed image vocabulary. In
548 *Proceedings of European Conference on Computer Vision*,
549 pages 97–112, 2002. 6
- 550 [4] Mark Everingham, Luc Van Gool, Christopher KI Williams,
551 John Winn, and Andrew Zisserman. The pascal visual object
552 classes (voc) challenge. *International Journal of Computer
553 Vision*, 88:303–338, 2010. 6
- 554 [5] Xiang Fang, Yuchong Hu, Pan Zhou, and Dapeng Wu. An-
555 imc: A soft approach for autoweighted noisy and incomplete
556 multiview clustering. *IEEE Transactions on Artificial Intel-
557 ligence*, 3(2):192–206, 2021. 1
- 558 [6] Xiang Fang, Daizong Liu, Pan Zhou, and Yuchong Hu.
559 Multi-modal cross-domain alignment network for video mo-
560 ment retrieval. *IEEE Transactions on Multimedia*, 25:7517–
561 7532, 2022. 1
- 562 [7] Marco Federici, Anjan Dutta, Patrick Forré, Nate Kush-
563 man, and Zeynep Akata. Learning robust representations via
564 multi-view information bottleneck. In *ICLR*, 2020. 2
- 565 [8] Michael Grubinger, Paul Clough, Henning Müller, and
566 Thomas Deselaers. The iapr tc-12 benchmark: A new eval-
567 uation resource for visual information systems. In *Internat-
568 ional Workshop onto Image*, pages 1–11, 2006. 6
- 569 [9] Renxiang Guan, Tianrui Liu, Wenxuan Tu, Chang Tang,
570 Wenhan Luo, and Xinwang Liu. Sampling enhanced con-
571 trastive multi-view remote sensing data clustering with long-
572 short range information mining. *IEEE Transactions on
573 Knowledge and Data Engineering*, pages 1–15, 2025. 1
- 574 [10] Pingting Hao, Kunpeng Liu, and Wanfu Gao. Anchor-guided
575 global view reconstruction for multi-view multi-label feature
576 selection. *Information Sciences*, 679:121124, 2024. 6
- 577 [11] Pingting Hao, Wanfu Gao, and Liang Hu. Embedded feature
578 fusion for multi-view multi-label feature selection. *Pattern
579 Recognition*, 157:110888, 2025. 1
- 580 [12] Mark J Huiskes and Michael S Lew. The mir flickr retrieval
581 evaluation. In *Proceedings of ACM International Conference
582 on Multimedia Information Retrieval*, pages 39–43, 2008. 6
- 583 [13] Yinggang Li, Xiangrong Tong, and Zhongming Lv. Multi-
584 dimensional requirements for reinforcement recommenda-
585 tion reasoning. *Applied Intelligence*, 55(6):1–16, 2025. 2
- 586 [14] Chengliang Liu, Jie Wen, Yabo Liu, Chao Huang, Zhihao
587 Wu, Xiaoling Luo, and Yong Xu. Masked two-channel de-
588 coupling framework for incomplete multi-view weak multi-
589 label learning. *Advances in Neural Information Processing
590 Systems*, 36:32387–32400, 2023. 2, 6
- 591 [15] Chengliang Liu, Jie Wen, Xiaoling Luo, Chao Huang, Zhi-
592 hao Wu, and Yong Xu. Dicnet: Deep instance-level con-
593 trastive network for double incomplete multi-view multi-
594 label classification. In *Proceedings of the AAAI conference
595 on artificial intelligence*, pages 8807–8815, 2023. 2, 6
- 596 [16] Chengliang Liu, Jie Wen, Xiaoling Luo, and Yong Xu. In-
597 complete multi-view multi-label learning via label-guided
598 masked view-and category-aware transformers. In *Pro-
599 ceedings of the AAAI conference on artificial intelligence*, pages
600 8816–8824, 2023. 2, 6
- 601 [17] Chengliang Liu, Jinlong Jia, Jie Wen, Yabo Liu, Xiaoling
602 Luo, Chao Huang, and Yong Xu. Attention-induced em-
603 bedding imputation for incomplete multi-view partial multi-
604 label classification. In *Proceedings of the AAAI Conference
605 on Artificial Intelligence*, pages 13864–13872, 2024. 2, 6
- 606 [18] Chengliang Liu, Gehui Xu, Jie Wen, Yabo Liu, Chao Huang,
607 and Yong Xu. Partial multi-view multi-label classification
608 via semantic invariance learning and prototype modeling. In
609 *Forty-first International Conference on Machine Learning*,
610 2024. 2, 6
- 611 [19] Chengliang Liu, Jie Wen, Yong Xu, Bob Zhang, Liqiang Nie,
612 and Min Zhang. Reliable representation learning for incom-
613 plete multi-view missing multi-label classification. *IEEE
614 Transactions on Pattern Analysis and Machine Intelligence*,
615 pages 1–17, 2025. 1
- 616 [20] Meng Liu, Yong Luo, Dacheng Tao, Chao Xu, and Yonggang
617 Wen. Low-rank multi-view learning in matrix completion for
618 multi-label image classification. In *Proceedings of the AAAI
619 conference on artificial intelligence*, 2015. 1
- 620 [21] Yuanyuan Liu, Jiyao Peng, Wei Dai, Jiabei Zeng, and
621 Shiguang Shan. Joint spatial and scale attention network for
622 multi-view facial expression recognition. *Pattern Recog-
623 nition*, 139:109496, 2023. 3
- 624 [22] Zhongchen Ma and Songcan Chen. Expand globally, shrink
625 locally: Discriminant multi-label learning with missing la-
626 bels. *Pattern Recognition*, 111:107675, 2021. 6
- 627 [23] Yalan Qin, Xinpeng Zhang, Shui Yu, and Guorui Feng. A
628 survey on representation learning for multi-view data. *Neu-
629 ral Networks*, 181:106842, 2025. 1
- 630 [24] Chongjie Si, Yuheng Jia, Ran Wang, Min-Ling Zhang,
631 Yanghe Feng, and Chongxiao Qu. Multi-label classifica-
632 tion with high-rank and high-order label correlations. *IEEE
633 Transactions on Knowledge and Data Engineering*, 36(8):
634 4076–4088, 2023. 2
- 635 [25] Gencer Sumbul, Marcela Charfuelan, Begüm Demir, and
636 Volker Markl. Bigearthnet: A large-scale benchmark archive
637 for remote sensing image understanding. In *IGARSS 2019-
638 2019 IEEE international geoscience and remote sensing
639 symposium*, pages 5901–5904. IEEE, 2019. 1
- 640 [26] Zhaobin Sun, Nannan Wu, Junjie Shi, Li Yu, Kwang-Ting
641 Cheng, and Zengqiang Yan. Fedmlp: Federated multi-label
642 medical image classification under task heterogeneity. In *In-
643 ternational Conference on Medical Image Computing and
644 Computer-Assisted Intervention*, pages 394–404. Springer,
645 2024. 1
- 646 [27] Qiaoyu Tan, Guoxian Yu, Carlotta Domeniconi, Jun Wang,
647 and Zili Zhang. Incomplete multi-view weak-label learning.
648 In *Ijcai*, pages 2703–2709, 2018. 6
- 649 [28] Zhihui Tian, John Upchurch, G Austin Simon, José Dubeux,
650 Alina Zare, Chang Zhao, and Joel B Harley. Quantifying

- 651 heterogeneous ecosystem services with multi-label soft clas-
652 sification. In *IGARSS 2024-2024 IEEE International Geo-*
653 *science and Remote Sensing Symposium*, pages 427–431.
654 IEEE, 2024. 1
- 655 [29] Jie Wen, Zheng Zhang, Lunke Fei, Bob Zhang, Yong Xu,
656 Zhao Zhang, and Jinxing Li. A survey on incomplete mul-
657 tiview clustering. *IEEE Transactions on Systems, Man, and*
658 *Cybernetics: Systems*, 53(2):1136–1149, 2022. 1
- 659 [30] Jie Wen, Chengliang Liu, Shijie Deng, Yicheng Liu, Lunke
660 Fei, Ke Yan, and Yong Xu. Deep double incomplete multi-
661 view multi-label learning with incomplete labels and missing
662 views. *IEEE transactions on neural networks and learning*
663 *systems*, 2023. 2, 6
- 664 [31] Changqing Zhang, Ziwei Yu, Qinghua Hu, Pengfei Zhu, Xin-
665 wang Liu, and Xiaobo Wang. Latent semantic aware multi-
666 view multi-label classification. In *Proceedings of the AAAI*
667 *conference on artificial intelligence*, 2018. 1
- 668 [32] Dawei Zhao, Qingwei Gao, Yixiang Lu, and Dong Sun.
669 Non-aligned multi-view multi-label classification via learn-
670 ing view-specific labels. *IEEE Transactions on Multimedia*,
671 25:7235–7247, 2022. 6
- 672 [33] Lihua Zhou, Guowang Du, Kevin Lue, Lizheng Wang, and
673 Jingwei Du. A survey and an empirical evaluation of multi-
674 view clustering approaches. *ACM Computing Surveys*, 56
675 (7):1–38, 2024. 2