

View-Category Interactive Sharing Transformer for Incomplete Multi-View Multi-Label Learning

Supplementary Material

6. Algorithm

In this section, we provide the additional explanation of our method. In Algorithm 1, the detailed training process is presented.

Algorithm 1 Training process of VIST

Input: Incomplete multi-view data $\{\mathbf{X}_v\}_{v=1}^V$, label matrix \mathbf{Y} , missing-view indicator matrix \mathbf{W} , missing-label indicator matrix \mathbf{U} .

Output: Final prediction $\hat{\mathbf{y}}$.

- 1 Initialize the parameters of VIST and set hyperparameters α, β, γ, k and training epochs E_{tr} .
 - 2 $t = 0$.
 - 3 **while** $t < E_{tr}$ **do**
 - 4 **for** $v = 1 : V$ **do**
 - 5 Embed the multi-view data \mathbf{X}_v to get embedding vector \mathbf{e}_v .
 - 6 **end**
 - 7 Stack the embedding vectors to obtain origin multi-view embedding \mathbf{E} .
 - 8 Generate the missing views by sampling from distribution $\mathcal{N}(\mu, \Sigma)$ with Eq. (9) and Eq. (10).
 - 9 Calculate the loss functions of contrastive learning \mathcal{L}_{pos} by Eq. (12) and \mathcal{L}_{neg} by Eq. (13).
 - 10 Get multi-view embedding $\tilde{\mathbf{E}}$ by Eq. (1), advanced multi-view embedding $\bar{\mathbf{E}}$ by Eq. (7) and calculate \mathcal{L}_{aux} using Eq. (17).
 - 11 Calculate the loss functions for view-category consistency guided embedding enhancement \mathcal{L}_c using Eq. (18).
 - 12 Get the final prediction $\hat{\mathbf{y}}$ by Eq. (19).
 - 13 Calculate multi-label classification loss \mathcal{L}_m using Eq. (20) and masked asymmetric loss \mathcal{L}_a using Eq. (22).
 - 14 Calculate the overall loss function \mathcal{L}_o using Eq. (23).
 - 15 Update the parameters.
 - 16 $t = t + 1$.
 - 17 **end**
-

7. Additional Visualizations

We carry out a series of experiments to visualize the generated complete multi-view data in the Corel5k datasets with 50% missing instances. Due to the category correlation within multi-label datasets, we opt to select five categories with semantically distant labels for the data visualization

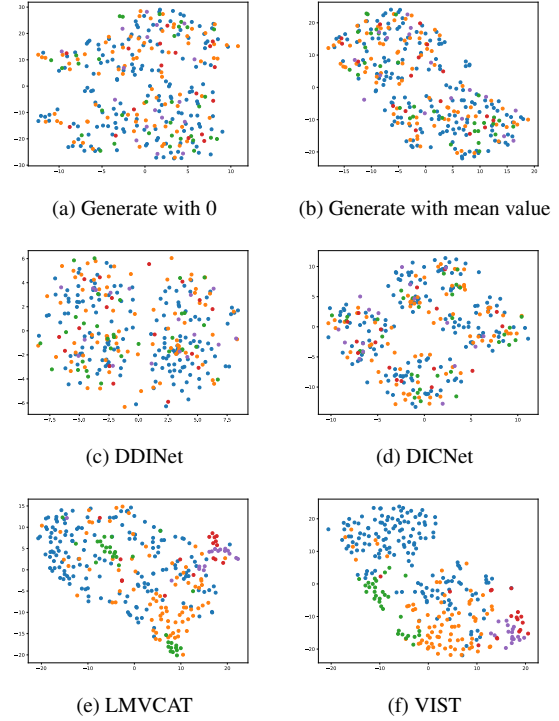


Figure 6. The results of different value of hyperparameter α and k on different datasets. β and γ are respectively set to default values of 0.1 and 0.1.

using t-SNE. These categories are water, beach, coyote, cars and hats. The results, as depicted in Fig. 6, indicate a pronounced enhancement in the distinction of margins between different classes. This improvement is attributable to the VIST’s capability to effectively explore and utilize the correlations existing among various views and categories.

Furthermore, we also present visualizations comparing the performance of our method with the current state-of-the-art approaches on datasets Corel5k and Espgame, as shown in Fig. 7 and Fig. 8. It is observable that our model achieves superior performance in these instances, a result that is directly attributable to the effective interaction between the views and categories.

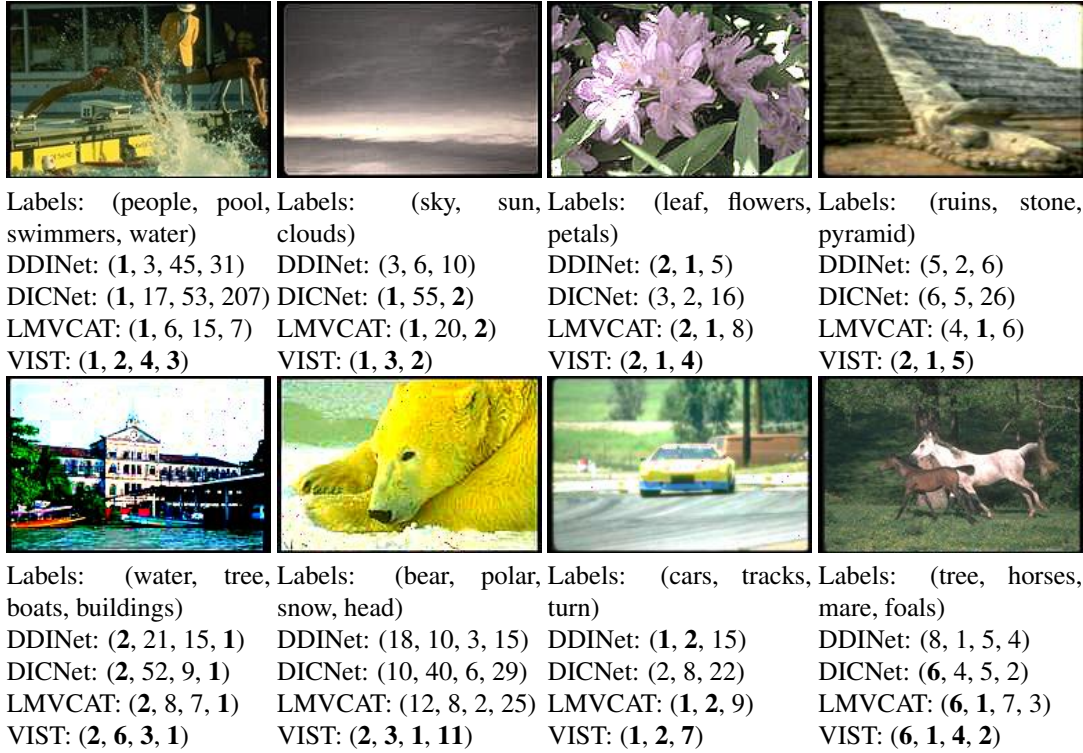


Figure 7. Visualization of prediction results on Dataset Core5k. The numbers enclosed in parentheses represent the ranking assigned by the model to the likelihood of the image belonging to each respective category, the same below.

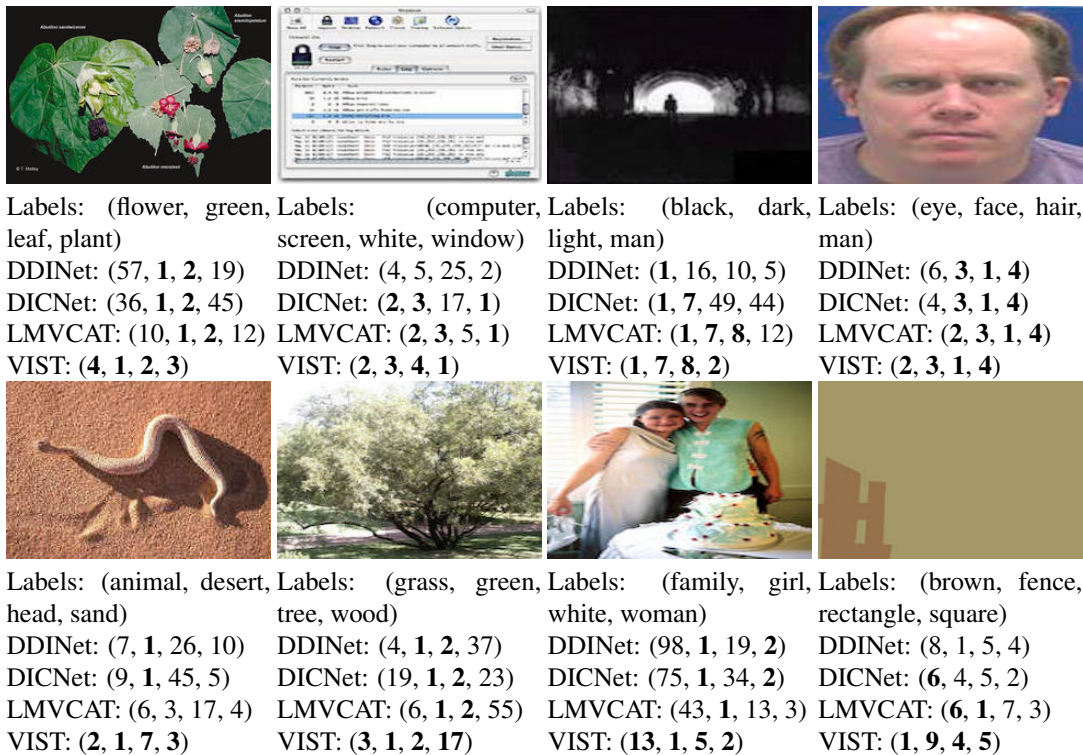


Figure 8. Visualization of prediction results on Dataset Espgame.