

Robult: Leveraging Redundancy and Modality-Specific Features for Robust Multimodal Learning

Duy A. Nguyen^{1,3}, Abhi Kamboj², Minh N. Do^{2,3}

¹Siebel School of Computing and Data Science, UIUC, US

²Department of Electrical and Computer Engineering, UIUC, US

³VinUni-Illinois Smart Health Center, VinUniversity, Vietnam

{duyan2, akamboj2, minhdo}@illinois.edu

Abstract

Addressing missing modalities and limited labeled data is crucial for advancing robust multimodal learning. We propose **Robult**, a scalable framework designed to mitigate these challenges by preserving modality-specific information and leveraging redundancy through a novel information-theoretic approach. Robult optimizes two core objectives: (1) a soft Positive-Unlabeled (PU) contrastive loss that maximizes task-relevant feature alignment while effectively utilizing limited labeled data in semi-supervised settings, and (2) a latent reconstruction loss that ensures unique modality-specific information is retained. These strategies, embedded within a modular design, enhance performance across various downstream tasks and ensure resilience to incomplete modalities during inference. Experimental results across diverse datasets validate that Robult achieves superior performance over existing approaches in both semi-supervised learning and missing modality contexts. Furthermore, its lightweight design promotes scalability and seamless integration with existing architectures, making it suitable for real-world multimodal applications.

1 Introduction

Motivation: In the Big Data era, multimodal learning significantly improves data exploitation, outperforming single-modality approaches [Huang *et al.*, 2021]. However, most existing methods [Daunhawer *et al.*, 2023; Peng *et al.*, 2022] operate under idealized assumptions: fully labeled training datasets and consistently available modalities during evaluation. In practice, the challenges of missing modalities and semi-supervised learning often coexist, yet current research typically addresses these problems in isolation. For example, missing modalities may arise when autonomous vehicles lose sensor inputs due to environmental obstructions or when medical diagnostics in resource-limited settings lack access to all imaging modalities. Simultaneously, the scarcity of labeled data across domains remains a critical bottleneck, particularly in multimodal contexts where individual modalities often require specialized labeling.

Addressing these challenges independently limits the adaptability of multimodal systems and fails to capture the complexities of real-world deployments. This work uniquely addresses these dual challenges by integrating strategies to handle incomplete modalities and leverage unlabeled data simultaneously. By doing so, our approach enhances both flexibility and applicability, enabling robust performance in diverse and imperfect scenarios. Unlike current literature, which rarely addresses both issues concurrently, this work bridges a critical gap with an innovative and unified solution.

Existing literature: One of the primary challenges in multimodal learning is handling corrupted or missing modalities. Existing methods to address missing modalities typically fall into two categories:

1. **Generative approaches**, such as VAE-based models [Wu and Goodman, 2018], which reconstruct missing modalities. While recent advancements [Feichtenhofer *et al.*, 2022; Woo *et al.*, 2023] demonstrate promising performance, these methods often depend on specific architectures, limiting their flexibility.
2. **Transfer learning methods**, which align latent spaces for cross-modal knowledge transfer [Ma *et al.*, 2022; Lee and Van der Schaar, 2021; Wang *et al.*, 2020]. These methods focus on adaptable training strategies [Chen *et al.*, 2023] but often lack a strong theoretical foundation and are primarily guided by empirical intuition.

Simultaneously, the need for semi-supervised learning arises from practical challenges in labeling raw data, especially in domains where annotations are scarce or labor-intensive to obtain. This challenge is amplified in multimodal learning, as each modality may require distinct expertise for labeling. For example, tasks such as object segmentation across video and lidar data in autonomous driving [Zhang *et al.*, 2022b] or medical segmentation across imaging modalities [Acosta *et al.*, 2022] demand diverse and often non-standardized labeling procedures. Recent advancements in semi-supervised learning, such as knowledge distillation [Su *et al.*, 2021] and pseudo-labeling [Aberdam *et al.*, 2022], have shown promise, but these techniques are often designed for specific applications and struggle to generalize across varied multimodal settings.

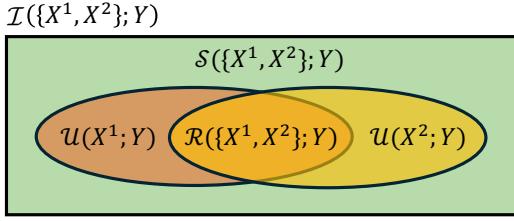


Figure 1. Partial Information Decomposition for 2 input modalities with target variable.

A detailed discussion of existing literature in both areas is provided in Appendix A.

Proposed approach: The dual challenges of missing modalities and semi-supervised learning remain critical open problems in multimodal research. This work advances transfer learning methods for missing modalities (category 2) that retain architectural flexibility and dataset compatibility, while introducing a novel design to reduce dependence on labeled data, ensuring robustness and adaptability in real-world scenarios.

Using Partial Information Decomposition [Williams and Beer, 2010], the mutual information provided by an input X with M modalities (X^1, \dots, X^M) for a given task Y can be decomposed into:

$$\begin{aligned} I(\{X^1, \dots, X^M\}; Y) &= R(\{X^1, \dots, X^M\}; Y) \\ &+ \sum_{i=1}^M U(X^i; Y) + S(\{X^1, \dots, X^M\}; Y), \end{aligned} \quad (1)$$

where R represents redundancy (shared task-relevant information among M modalities), U denotes unique information specific to the i^{th} modality, and S quantifies synergy, the additional knowledge generated through interactions among modalities (Figure 1 illustrates PID with 2 modalities).

In an ideal setting with access to all modalities, a fusion technique would efficiently capture R and S to optimize predictions for Y . However, in real-world scenarios where some modalities are unavailable, replicating S becomes challenging. For example, with two modalities, $Y = \text{linear}(X^1, X^2) = \text{non_linear}(X^1)$ demonstrates that synergy (S) improves predictions when both X^1 and X^2 are available. Access to X^1 alone makes the relationship more complex and harder to model with deep networks.

Approaches like knowledge distillation [Chen *et al.*, 2023] and contrastive learning [Radford *et al.*, 2021] aim to address missing modalities by mimicking the representations produced by fused modalities in their absence. From an information-theoretic perspective, these methods focus on replicating redundant information (R) through latent-space alignment. Building on this foundation, we explicitly introduce a mutual information maximization objective (Objective 2.1) to align unimodal and fused representations. This alignment ensures efficient knowledge transfer while minimizing reliance on labeled data using a novel soft-positive

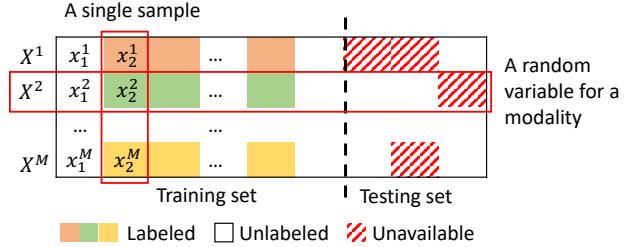


Figure 2. Training/evaluation datasets under investigation.

pseudo-labeling mechanism that accounts for pseudo-label uncertainty.

While alignment effectively captures R , it can unintentionally diminish unique information (U) from individual modalities. This loss is particularly detrimental in semi-supervised settings or when other modalities are unavailable. To address this, we propose preserving modality-specific information (U^i) via Objective 2.2. This is achieved using a simple yet effective reconstruction procedure in the latent space, universally applicable across modalities. Our results and ablations (Tables 1 and 3) demonstrate how retaining U^i improves performance.

Together, Objectives 2.1 and 2.2 form the foundation of our semi-supervised multimodal learning method, **Robust Multimodal Pipeline (Robult)**. Robult effectively balances redundancy alignment and unique information preservation, ensuring accuracy and robustness in scenarios with missing modalities. Empirical results (Section 3) and theoretical underpinnings highlight Robult's superior performance compared to existing methods, demonstrating its adaptability to diverse real-world settings.

Contributions. Our primary contributions are:

- Jointly addressing the dual challenges of missing modalities and semi-supervised learning.
- Framing two objectives under an information-theoretic perspective and deriving novel loss functions to achieve these goals.
- Introducing a soft Positive-Unlabeled contrastive loss that efficiently utilizes limited labeled data through selective weighting of potential positives.

2 Methodology

Training setting: We consider a training scenario where each sample includes all modalities, but only some samples have corresponding labels (Figure 2). Formally, let the training dataset be $\mathcal{D}_{\text{train}} = \{(x_1, y_1), \dots, (x_k, y_k), x_{k+1}, \dots, x_n\}$, where each data point $x_j = (x_j^1, \dots, x_j^M)$ contains M modalities. The dataset consists of n samples, of which k are labeled ($0 \leq k < n$).

Evaluation setting: To mimic real-world deployment, the evaluation dataset $\mathcal{X}_{\text{test}} = \{x_1, \dots, x_m\}$ consists of samples with potential missing modalities - e.g. $x_j = (x_j^a)^*$; $\forall a \in$

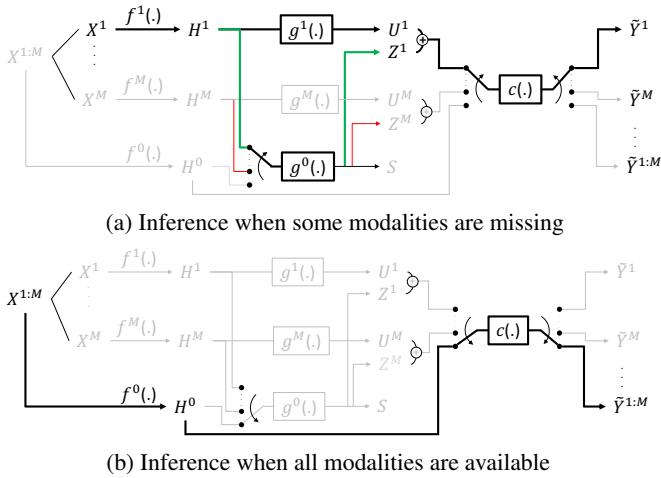


Figure 3. (3a) For each available modality $i \in \{1, \dots, M\}$, $f^i(\cdot)$ extracts the latent representation h_j^i . The modules $g^i(\cdot)$ and $g^0(\cdot)$ then extract the unique information U_j^i and redundant representation Z_j^i , respectively. These are combined and processed by the shared module $c(\cdot)$ to predict \tilde{y}_j^i . Late fusion strategies aggregate these outputs into the final prediction \tilde{y}_j . (3b) When all modalities $x_j^{1:M}$ are available, $f^0(\cdot)$ extract a fused latent vector H_j^0 . This vector is directly passed through $c(\cdot)$ to generate the final output $\tilde{y}_j^{1:M}$.

$\mathcal{A}_j^M \subseteq \{1, \dots, M\}$. Here, \mathcal{A}_j^M denote the indices of all available modalities for sample x_j (Figure 2).

Notation: For clarity, in the theoretical framework, we denote the input random variables corresponding to the i^{th} modality as X^i (Figure 2). We would refer back to the observation-level notation x_j^i whenever needed.

Figure 3 illustrates our versatile multimodal pipeline named Robult. Robult consists of M modality-specific branches and a fusion branch indexed with zero. For each input variable X^i , it is first projected into a shared latent space:

$$H^i = f^i(X^i), \quad i \in \{1, \dots, M\}.$$

In parallel, a fusion network $f^0(\cdot)$ processes all input modalities jointly to produce a fused latent variable (also reside on the shared latent space):

$$H^0 = f^0(X^{1:M}),$$

where index 0 denotes the joint latent variable using all modalities. This latent-space projection provides two key benefits:

- **Generalization:** Robult supports different combination of modalities and their preferred projection methods/architectures.
- **Efficient Fusion:** Various strategies can be adopted before Robult’s main logic to efficiently attain fused representations.

Each modality-specific module $g^i(\cdot)$ extracts unique information of corresponding modality:

$$U^i = g^i(H^i), \quad i \in \{1, \dots, M\},$$

while the shared module $g^0(\cdot)$ effectively captures the redundant information by processing individual or joint latent variables:

$$Z^i = g^0(H^i), S = g^0(H^0).$$

Objectives: As discussed in Section 1, we draw inspiration from the second category of literature on missing modalities, which primarily leverages knowledge distillation and contrastive loss techniques [Poklukar *et al.*, 2022b]. These methods aim to align unimodal (student) representations, such as Z^i , with the fused (teacher) representation S , enabling unimodal representations to efficiently replicate the redundant information encapsulated in the fused representation. Although this redundant information exists within each modality, extracting it directly from a single modality can be significantly more complex without access to all modalities. For example, encoding object shape from an image may be challenging, whereas a textual description can explicitly provide the same information. Mimicking the fused representation offers a shortcut, simplifying the extraction of redundant information within unimodal representations.

However, these approaches fail in the semi-supervised training setting, owing to 2 emerging challenges:

- (1) The *over reliance on label signals* (e.g. classification clusters) in the alignment process.
- (2) the *diminishing of unique information* contained in different modality after alignment process.

Regarding the former challenge, we directly address the label need of Contrastive Learning with Label-level sampling [Zhang *et al.*, 2022a] by a novel soft Positive-Unlabelled (PU) contrastive loss, together with an adaptive weighting strategy, detailed about which is covered in Section 2.1. Theoretically, this PU contrastive loss corresponds to a mutual information maximization problem between the desired fused representation S and the learned unimodal representation projected into that same latent space Z^i ’s. Objective 2.1 of our method can be expressed as follow:

Objective 2.1. Aligning S and Z^i by maximizing the mutual information $\mathcal{I}(S, Z^i) \quad (i = \overline{1, M})$.

For the latter challenge, we observe that attempting to only align modalities during training diminishes the unique information provided by each modality, thus the model is losing information that could help inform it when only that specific modality is available. Therefore, a model should ideally benefit by maintaining the redundancy while also explicitly preserving each modality’s unique information. This claim is later supported by experimental results and ablations (Section 3 - Table 1). To address the challenge of vanishing modality-specific information from multimodal alignment during training, we emphasize a disentanglement strategy that preserves unique information while still facilitating the redundancy learning process of Objective 2.1. Robult integrates a set of modules $g^i(H^i)$, where $i = 1, \dots, M$, to produce unique representations U^i for each modality. We aim to preserve the unique information for each modality via the learning of U^i with Objective 2.2 as follows:

Objective 2.2. Learning U^i by minimizing the conditional entropy $\mathcal{H}(H^i|Z^i, U^i)$ ($i = 1, M$).

During training, all branches are executed and three loss functions update the learned modules. The soft positive unlabeled loss, \mathcal{L}_{PU} , maximizes knowledge extraction from the few labeled samples in a batch (Subsection 2.1). The reconstruction loss, \mathcal{L}_{rec} forces each branch to extract unique modality-specific information (Subsection 2.2). The task-specific supervised loss, \mathcal{L}_{sup} is used on the labeled samples across all modules to learn label information (Subsection 2.3).

2.1 Maximizing Mutual Information with Soft Positive-Unlabeled Contrastive Learning

To address the impact of missing modalities, we aim to learn redundant information by aligning the fused latent variable S with unimodal representations Z^i , as formalized in Objective 2.1. This is achieved by maximizing their mutual information $\mathcal{I}(S, Z^i)$. However, direct computation of this quantity is infeasible without access to the joint distribution p_{S, Z^i} or the marginal distributions p_S and p_{Z^i} . Instead, we derive and optimize a lower bound for this mutual information.

Lower Bound Derivation. Let F be a binary random variable indicating whether a pair (s_j, z_k^i) is sampled from the joint distribution p_{S, Z^i} ($F = 1$) or from the product of marginal distributions $p_S \otimes p_{Z^i}$ ($F = 0$). Then, a lower bound for $\mathcal{I}(S, Z^i)$ is expressed as:

$$\begin{aligned} \mathcal{I}(S, Z^i) &\geq -\mathbb{E}_{p_{S, Z^i}} \log v(S, Z^i) \\ &= -\mathbb{E}_{p(S, Z^i | F=1)} \log v(S, Z^i) \end{aligned} \quad (2)$$

where $v(S, Z^i)$ is a non-parametric approximation of $p(F = 1 | S, Z^i)$. For a sampled pair (s_j, z_k^i) in a batch of B samples, where s_j is the fused representation for sample j and z_k^i is the i^{th} modality-specific representation for sample k , $v(s_j, z_k^i)$ is defined as:

$$v(s_j, z_k^i) = \frac{\phi(s_j, z_k^i)}{\sum_h^B \phi(s_j, z_h^i)}; \\ \text{where } \phi(s_j, z_k^i) = \exp(\langle s_j; z_k^i \rangle / \tau).$$

The derivation of Result 2 is detailed in Appendix B.2.

Challenges with Sampling. The lower bound in equation 2 relies on expectations under $\mathbb{E}_{p_{S, Z^i}}$, a key source of deviation in existing studies, which presents challenges in practice. Two common sampling strategies include:

1. **Instance-level sampling**, which considers only intra-sample pairs (s_j, z_j^i) within q mini-batch [Radford *et al.*, 2021].
2. **Label-level sampling**, which uses label information to sample inter-sample pairs (s_j, z_k^i) with the same labels, e.g. $y_j = y_k$ [Zhang *et al.*, 2022a].

The first approach risks introducing false negatives, while the second requires fully labeled data, limiting its applicability in semi-supervised settings.

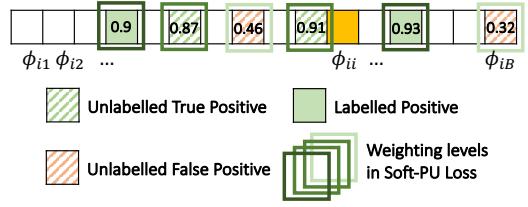


Figure 4. Soft-PU Loss mechanism. Unlabeled positive pairs are identified using soft labels from Robult’s classifier. These pairs are re-weighted based on their proximity and the mean proximity of true labeled positive pairs to mitigate false positives.

Soft Positive-Unlabeled (PU) Contrastive Loss. To overcome these limitations, we propose a novel soft Positive-Unlabeled (PU) contrastive loss with adaptive weighting. Let L indicate whether a sampled pair (s_j, z_k^i) is labeled ($L = 1$) or not ($L = 0$) (where $L = 1$ if the label information of both samples i and k is known, and $L = 0$ otherwise). The lower bound in equation 2 can then be decomposed as:

$$\begin{aligned} -\mathbb{E}_{p_{S, Z^i}} \log v(S, Z^i) &= -\mathbb{E}_{p(S, Z^i | F=1)} \log v(S, Z^i) \\ &= -p(L=1)\mathbb{E}_{p(S, Z^i | F=1, L=1)} \log v(S, Z^i) \\ &\quad - p(L=0)\mathbb{E}_{p(S, Z^i | F=1, L=0)} \log v(S, Z^i) \end{aligned} \quad (3)$$

We formulate the two terms in equation 3 as separate loss components: \mathcal{L}_{lb} for labeled data and \mathcal{L}_{ulb} for the unlabeled data.

For labeled samples, let $B_{1,1}^j$ denote the index set of inputs in the batch that share the same class as sample j ($F = 1$) and are labeled ($L = 1$), i.e. $k \in B_{1,1}^j \iff (s_j, z_k^i) \sim p(S, Z^i | F = 1, L = 1)$. The labeled loss is then defined as a NT-Xent-like contrastive loss [Chen *et al.*, 2020]:

$$\mathcal{L}_{lb} = -\frac{1}{M} \sum_{i=1}^M \sum_{j \in B} -\frac{1}{\|B_{1,1}^j\|} \sum_{k \in B_{1,1}^j} \log v(s_j, z_k^i). \quad (4)$$

For unlabeled samples, directly sampling from $p(S, Z^i | F = 1, L = 0)$ is challenging due to the absence of label information. To address this, we propose leveraging the output of the Robult classifier to generate soft labels, regularized by adaptive weights calculated dynamically within each mini-batch. This approach effectively balances the influence of soft-labeled pairs and mitigates the limitations of traditional pseudo-labeling methods [Aberdam *et al.*, 2022].

A key challenge during initial training stages is the instability of the Robult classifier, which may produce unreliable outputs and hinder effective filtering of false positives, as illustrated in Figure 4. To counteract this, we adjust the contribution of soft-labeled pairs to the loss function, ensuring robust training. For each anchor sample s_j within a mini-batch, there are labeled positive partners or, in unsupervised scenarios, unimodal representations z_j^i . The average proximity of

these labeled partners to s_j serves as a reference for determining “positive” proximity thresholds. Unlabeled positive pairs identified by the Robult classifier are expected to exhibit proximity close to this reference mean. To refine the loss computation, we increase the weight of pairs that align closely with the reference mean and decrease the influence of outliers. This weighting mechanism, implemented using an RBF kernel, allows precise adjustment of each couplet’s contribution based on its proximity. By dynamically adapting weights, our method effectively reduces the impact of potential false positives and enhances the overall performance of the final loss function, as demonstrated in Figure 4.

$$w_{jk}^i = RBF(\phi_j^i, \phi(s_j, z_k^i)); \\ \text{where } \phi_j^i = \text{mean} \left\{ \phi(s_j, z_k^i) | k \in B_{1,1}^j \right\}. \quad (5)$$

Let $B_{1,0}^j$ denote the index set of inputs in the batch that share the same class as anchor j ($F = 1$) but are not labeled ($L = 0$), i.e. $k \in B_{1,0}^j \iff (s_j, z_k^i) \sim p(S, Z^i | F = 1, L = 0)$. The unlabeled loss is given by:

$$\mathcal{L}_{ulb} = -\frac{1}{M} \sum_{i=1}^M \sum_{j \in B} -\frac{1}{||B_{1,0}^j||} \sum_{k \in B_{1,0}^j} w_{jk}^i \log v(s_j, z_k^i). \quad (6)$$

The complete soft Positive-Unlabeled (PU) loss is the sum of the labeled and unlabeled components:

$$\mathcal{L}_{PU} = \mathcal{L}_{lb} + \mathcal{L}_{ulb} \quad (7)$$

2.2 Minimizing Conditional Entropy with Latent Reconstruction Error

This section details the procedure for achieving Objective 2.2, which focuses on preserving unique information U^i . Let p_{U^i, Z^i} denote the joint distribution of U^i and Z^i , where $(u_j^i, z_j^i) \sim p_{U^i, Z^i}$ are derived from the corresponding instance h_j^i . Inspired by [Chen *et al.*, 2016], the conditional entropy $\mathcal{H}(H^i | Z^i, U^i)$ is expressed as:

$$\mathcal{H}(H^i | U^i, Z^i) = -\mathbb{E}_{p_{U^i, Z^i}} \left[\mathbb{E}_{p_{H^i | U^i, Z^i}} [\log p(H^i | U^i, Z^i)] \right] \quad (8)$$

Upper Bound Derivation. Since directly computing $p(H^i | U^i, Z^i)$ is challenging, we approximate it using a distribution $q(H^i | U^i, Z^i)$. Substituting q into Eq. 8, we derive:

$$\begin{aligned} \mathcal{H}(H^i | Z^i, U^i) &= -\mathbb{E}_{U^i, Z^i} \left[\mathbb{E}_{H^i | U^i, Z^i} [\log q(H^i | U^i, Z^i) \cdot \frac{p(H^i | U^i, Z^i)}{q(H^i | U^i, Z^i)}] \right] \\ &= -\mathbb{E}_{U^i, Z^i} \left[\mathbb{E}_{H^i | U^i, Z^i} [\log q(H^i | U^i, Z^i)] + \text{KL}(p || q) \right] \\ &\leq -\mathbb{E}_{U^i, Z^i} \left[\mathbb{E}_{H^i | U^i, Z^i} [\log q(H^i | U^i, Z^i)] \right] \end{aligned} \quad (9)$$

The last inequality arises because the KL divergence $\text{KL}(p || q)$ is non-negative. Thus, minimizing $\mathcal{H}(H^i | Z^i, U^i)$ reduces to minimizing its Evidence Lower Bound (ELBO)-like [Kingma and Welling, 2014] upper bound using the approximating distribution $q(H^i | U^i, Z^i)$.

Modeling $q(H^i | U^i, Z^i)$ with Latent Reconstruction. We model $q(H^i | U^i, Z^i)$ through a latent reconstruction procedure in the shared latent space. Specifically, we define a

reconstruction module $r^i(U^i, Z^i) = \tilde{H}^i$ where \tilde{H}^i approximates H^i . For each pair $(u_j^i, z_j^i) \sim p_{U^i, Z^i}$ generated from h_j^i , the module $r^i(\cdot)$ attempts to reconstruct \tilde{h}_j^i such that it closely resembles h_j^i . The reconstruction loss is formulated as:

$$\mathcal{L}_{rec} = \frac{1}{MB} \sum_{i=1}^M \sum_{j=1}^B 1 - \langle \tilde{h}_j^i, h_j^i \rangle^2, \quad (10)$$

where $\langle \cdot, \cdot \rangle$ denotes the $L2$ -normalized dot product operation, B is the size of mini-batch, and M is the number of modalities.

This reconstruction in latent space is computationally efficient and alleviates the complexity of directly reconstructing raw modality data. Moreover, it generalizes well across various modalities, enhancing the flexibility of Robult. The reconstruction loss \mathcal{L}_{rec} is back-propagated exclusively through the M unimodal branches. This design ensures that the unique information $\mathcal{U}(X^i; Y)$ is preserved through U^i , without interfering with the shared branch’s focus on learning redundant information, as discussed in Section 2.1.

2.3 Training strategy

The objectives outlined in Sections 2.1 and 2.2 are optimized through their respective loss functions. Due to the distinct purposes of these objectives, it is advantageous to learn them separately. Specifically, we apply \mathcal{L}_{rec} to guide the learning of $g^i(\cdot)$ ($i = 1, \dots, M$), while \mathcal{L}_{PU} drives the optimization of $f^i(\cdot)$, $f^0(\cdot)$, and $g^0(\cdot)$. To maximize the utility of labeled data, we incorporate an additional supervised loss \mathcal{L}_{sup} . This loss directs the learning process for the entire Robult network and adapts based on the task type: L_1 loss for regression tasks and cross-entropy \mathcal{L}_{ce} for classification tasks. A detailed training procedure, including loss formulations and implementation details, is provided in Appendix B.3.

3 Experimental Results

3.1 Datasets and Metrics

Dataset: We conduct experiments on the following datasets.

- **CMU-MOSI** [Zadeh *et al.*, 2016] & **CMU-MOSEI** [Zadeh *et al.*, 2018b]: Containing three modalities - text, audio, and visual, supporting sentiment analysis and emotion recognition tasks. Each video is labeled on a scale from -3 (negative) to 3 (positive) sentiment.
- **MM-IMDb** [Arevalo *et al.*, 2017]: Containing image and text modalities, serving genre classification task - which involves multi-label classification as a movie has several genres.
- **UPMC Food-101** [Wang *et al.*, 2015]: Containing two modalities, text and images, this dataset is a classification dataset consisting of 101 food categories.
- **Hateful Memes** [Kiela *et al.*, 2020]: Containing text and image modalities, this dataset aims at identifying hate

speech in memes. This dataset includes challenging examples that are similar to hateful ones but are actually harmless.

Metrics: For sentiment analysis related to CMU-MOSI and CMU-MOSEI datasets, we adopt mean absolute error (MAE), correlation (Cor), binary accuracy, and F1 score, following [Poklukar *et al.*, 2022b; Tsai *et al.*, 2018]. Here, binary categories determine positive sentiment scores (> 0) or negative ones (< 0). For the evaluation of the three remaining datasets, we adhere to the metrics specified in [Lee *et al.*, 2023b]. With the MM-IMDb dataset, the multi-label classification performance is assessed using F1-Macro. The classification accuracy is employed for the UPMC Food-101 dataset. Lastly, for Hateful Memes, the evaluation is based on the AU-ROC metric.

3.2 Baselines and Experimental Settings

Baselines: We incorporate several state-of-the-art approaches representing popular strategies into our comparative evaluation. Specifically, GMC [Poklukar *et al.*, 2022b] serves as a contrastive learning-based approach, ActionMAE [Woo *et al.*, 2023] represents a generation-based method, and we include a Transformer-based approach proposed in [Lee *et al.*, 2023b], referred to as Prompt-Trans for brevity. To ensure optimal reproducibility, we inherit the implementations of all baseline methods from their original code bases. Additionally, we implement unimodal frameworks (Unimodal) for each modality, trained in a supervised manner with available labels, to serve as our baseline comparison.

Implementation details: To ensure a fair comparison, we use similar encoder architectures for processing raw data modalities whenever possible. The unimodal baselines are designed with the same architectures as Robult, each with its own classifier. For Robult, positive samples for the soft P-U loss are determined after discretizing labels if needed. Specifically, in the cases of CMU-MOSI and CMU-MOSEI datasets, label information in the range of $[-3, 3]$ is quantified into 7 discrete categories $(-3, -2, \dots, 3)$. Additionally, for the multi-label dataset MM-IMDb, two samples are considered positive if they share all the same labels. Regarding Prompt-Trans, we only report its results for three datasets involving two modalities (MM-IMDb dataset, UPMC Food-101 dataset, and Hateful Memes dataset), as the extension to multiple modalities cannot be directly inferred from the original work [Lee *et al.*, 2023b].

Experimental details: The primary focus of our performance reporting is on two extreme scenarios: semi-supervised settings with only 5% labeled data and scenarios where only a single modality is presented during evaluation. All reported results are averaged over 3 different random seeds. In the semi-supervised setup, the newly created labeled portion is ensured to maintain the correct label ratio as the original training sets. Additional experiments extending these two settings to more modalities and a higher percent of labeled data are detailed in Appendix D.1 and D.2 respectively. Specific details on implementation settings relating to each dataset are provided in Appendix - C.2.

Metrics	CMU-MOSI				CMU-MOSEI			
	Unimodal	GMC	ActionMAE	Robult	Unimodal	GMC	ActionMAE	Robult
<i>Text Modality:</i>								
MAE (\downarrow)	1.41	1.407	1.476	1.397	0.81	0.815	1.115	0.784
Corr (\uparrow)	0.137	0.14	0.066	0.144	0.383	0.346	0.136	0.459
F1 (\uparrow)	0.551	0.559	0.535	0.578	0.717	0.716	0.614	0.739
Acc (\uparrow)	0.553	0.562	0.47	0.569	0.712	0.708	0.603	0.732
<i>Audio Modality:</i>								
MAE (\downarrow)	1.576	1.518	1.546	1.415	0.842	0.836	1.215	0.825
Corr (\uparrow)	0.041	-0.065	0.046	0.085	0.111	0.193	0.101	0.221
F1 (\uparrow)	0.512	0.457	0.508	0.539	0.618	0.642	0.634	0.679
Acc (\uparrow)	0.496	0.46	0.467	0.535	0.599	0.63	0.543	0.65
<i>Vision Modality:</i>								
MAE (\downarrow)	1.451	1.497	1.511	1.425	0.891	0.839	1.127	0.826
Corr (\uparrow)	0.044	-0.07	-0.03	0.086	0.163	0.2	0.104	0.201
F1 (\uparrow)	0.585	0.446	0.511	0.593	0.637	0.621	0.594	0.647
Acc (\uparrow)	0.425	0.449	0.514	0.522	0.624	0.62	0.561	0.632
<i>Full Modality:</i>								
MAE (\downarrow)	1.394	1.47	1.496	1.392	0.783	0.819	1.103	0.779
Corr (\uparrow)	0.186	0.101	-0.092	0.247	0.364	0.328	0.337	0.504
F1 (\uparrow)	0.597	0.497	0.553	0.657	0.73	0.693	0.694	0.744
Acc (\uparrow)	0.594	0.498	0.477	0.63	0.729	0.688	0.643	0.741

Table 1. Results on CMU-MOSI, CMU-MOSEI.

3.3 Main Quantitative Results

All results are shown in tables with the best outcomes in red and the second-best in blue.

Sentiment Analysis: The results for CMU-MOSI and CMU-MOSEI datasets are summarized in Table 1. For both datasets, Robult significantly outperforms all the compared methods, suggesting its effectiveness and consistency in semi-supervised and missing modality scenarios. Regarding CMU-MOSI, due to its smaller scale compared to CMU-MOSEI, the labeled portions are also smaller. This condition poses a challenge for existing baselines that heavily rely on label information. In contrast, Robult effectively addresses this challenge, demonstrating the ability to extract meaningful representations even with limited labeled data. On CMU-MOSEI, Robult consistently produces superior representations, achieving the best performances across all recorded metrics. Notably, Robult improves the correlation (Corr) between the predicted sentiment levels and ground truth by up to 19.8%, outperforming the second-best method, which is the unimodal for textual data.

Classification tasks: In Table 2, empirical results for three classification tasks show that Robult consistently outperforms existing approaches and baselines in most cases, except for one scenario on the Hateful Memes dataset with the full modality available, where Robult achieves comparable performance with Prompt-Trans [Lee *et al.*, 2023b]. Notably, the Hateful Memes dataset includes samples with “benign confounders”, negatively impacting performance when models rely solely on single modalities [Kiela *et al.*, 2020]. Leveraging the soft Positive-Unlabelled loss, Robult effectively addresses and mitigates performance gaps with either single modality inputs or the full ones. In addition, we calculate F1 macro scores for all methods on these three datasets in the unimodal and multimodal cases. We further visualize a Critical Difference Diagram [Demšar, 2006] in Figure 5. This diagram visually represents the performance among different machine learning algorithms across various datasets by displaying the mean performance ranks, with lower being better,

	Unimodal	Prompt-Trans	GMC	ActionMAE	Robult
<i>MM-IMDb - F1 Macro (\uparrow):</i>					
Text	0.24	0.198	0.296	0.055	0.321
Image	0.207	0.148	0.291	0.039	0.298
Full	0.196	0.268	0.307	0.171	0.332
<i>UPMC Food-101 - Accuracy (\uparrow):</i>					
Text	0.321	0.151	0.395	0.196	0.435
Image	0.296	0.111	0.382	0.132	0.415
Full	0.138	0.432	0.41	0.358	0.446
<i>Hateful Memes - AUROC (\uparrow):</i>					
Text	0.584	0.511	0.617	0.528	0.623
Image	0.524	0.475	0.528	0.508	0.596
Full	0.618	0.635	0.616	0.542	0.632

Table 2. Results on MM-IMDb, UPMC Food-101, Hateful Memes.

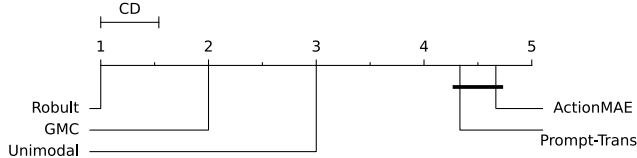


Figure 5. CD diagram showing the mean rank of each method on three datasets.

and connecting statistically indistinguishable groups (within 95% confidence level) with a thin horizontal bar, as per the Friedman hypothesis test. From the diagram, Robult exhibits a clear improvement gap compared to other state-of-the-art methods in average ranks, while ActionMAE and Prompt-Trans show no statistically significant difference in their performance.

3.4 Main ablation studies

Owing to space constraint, we provide here a key analysis of ablation study for different components of Robult. For a more comprehensive analysis, please refer to the additional experiments in Appendix D, which offer further insights into how architectural choices, Soft-PU loss, and weighting schemes influence Robult’s performance quantitatively and qualitatively.

We evaluate the impact of each loss component on Robult’s performance using Hateful Memes dataset, which mirror the semi-supervised and missing modalities conditions of our main experiments. This analysis involves testing variations of Robult with different ablations. (1) *Removal of \mathcal{L}_{sup}* - this setting utilizes available label information only in $\mathcal{L}_{(u)lb}$, so Robult can only produce latent representations. An additional Logistic Regressor is trained with these representations as its input, and this pipeline’s final scores are reported. (2) *Removal of \mathcal{L}_{rec}* - this setting discards \mathcal{L}_{rec} , corresponding to our Objective 2.2. (3) *Removal of \mathcal{L}_{lb}* - this setting makes the learning of Objective 2.1 rely only on \mathcal{L}_{ulb} . (4) *Removal of \mathcal{L}_{ulb}* - this setting associates Objective 2.1 exclusively with \mathcal{L}_{lb} . Table 3 summarizes the results of this ablation experiment. Overall, any ablation negatively impacts the performance of Robult. In particular, the absence of \mathcal{L}_{sup} signif-

Metrics	GMC	Robult	Robult (1)	Robult (2)	Robult (3)	Robult (4)
<i>CMU-MOSI - Text Modality:</i>						
MAE	1.407	1.397	1.589	1.511	1.443	1.429
Corr	0.14	0.144	0.047	0.101	0.051	0.123
F1	0.559	0.578	0.542	0.52	0.593	0.571
Acc	0.562	0.569	0.544	0.523	0.422	0.573
<i>CMU-MOSI - Audio Modality:</i>						
MAE	1.518	1.415	1.586	1.561	1.494	1.495
Corr	-0.065	0.085	0.023	0.005	0.046	0.085
F1	0.457	0.539	0.526	0.499	0.51	0.517
Acc	0.46	0.535	0.518	0.502	0.442	0.509
<i>CMU-MOSI - Vision Modality:</i>						
MAE	1.497	1.425	1.663	1.711	1.445	1.504
Corr	-0.07	0.086	0.025	-0.023	0.041	-0.066
F1	0.446	0.593	0.519	0.485	0.571	0.459
Acc	0.449	0.522	0.519	0.465	0.47	0.448
<i>CMU-MOSI - Full Modality:</i>						
MAE	1.47	1.392	1.588	1.434	1.411	1.459
Corr	0.101	0.247	0.071	0.239	0.166	0.229
F1	0.497	0.657	0.524	0.567	0.549	0.6
Acc	0.498	0.63	0.523	0.566	0.552	0.601
<i>Hateful Memes - Text Modality:</i>						
AUROC	0.617	0.623	0.528	0.59	0.605	0.576
Acc	0.581	0.59	0.535	0.556	0.571	0.562
<i>Hateful Memes - Image Modality:</i>						
AUROC	0.528	0.596	0.518	0.582	0.588	0.566
Acc	0.551	0.562	0.524	0.551	0.526	0.539
<i>Hateful Memes - Full Modality:</i>						
AUROC	0.616	0.632	0.538	0.618	0.634	0.582
Acc	0.532	0.595	0.542	0.55	0.554	0.552

Table 3. Ablation analysis on CMU-MOSI and Hateful Memes datasets for Robult.

icantly worsens the performance, as there is no loss guiding the learning of Robult’s classifier, which is crucial for generating soft label information consumed by the soft Positive-Unlabeled loss \mathcal{L}_{ulb} . Consequently, this ablation adversely affects two loss components, explaining the poorest result among all variations. The removal of \mathcal{L}_{rec} particularly harms the performance with unimodal inputs, aligning with the motivation for Objective 2.2, as the unique information U^i is no longer preserved. In two remaining cases, both ablations diminish Robult’s overall performance, indicating their equal contribution to achieving Objective 2.1.

4 Contributions & Limitations

Contributions: Our Robult pipeline leverages limited label data through a soft Positive-Unlabelled (PU) loss and latent reconstruction loss, enhancing modality interactions and preserving unimodal data integrity. It supports various modality types and quantities, scales linearly with modalities, and functions independently of specific architectures. This flexibility facilitates integration with existing DL frameworks, advancing multimodal learning in practical settings.

Limitations. Robult’s design presumes that the proximity of positive couples follows a Gaussian distribution, a method proven empirically but not theoretically. Future work should seek theoretical validation for this assumption. Moreover, with our setting, the potential of labeled data in scenarios with missing modalities in training remains untapped. Exploring these cases could further improve Robult’s effectiveness in complex real-world applications.

Acknowledgements

The work of Duy A. Nguyen was supported in part by a PhD fellowship from the VinUni-Illinois Smart Health Center, VinUniversity, Hanoi, Vietnam.

References

- [Aberdam *et al.*, 2022] Aviad Aberdam, Roy Ganz, Shai Mazor, and Ron Litman. Multimodal semi-supervised learning for text recognition. *arXiv preprint arXiv:2205.03873*, 2022.
- [Acosta *et al.*, 2022] Julián N Acosta, Guido J Falcone, Pranav Rajpurkar, and Eric J Topol. Multimodal biomedical ai. *Nature Medicine*, 28(9):1773–1784, 2022.
- [Arevalo *et al.*, 2017] John Arevalo, Thamar Solorio, Manuel Montes-y Gómez, and Fabio A González. Gated multimodal units for information fusion. *arXiv preprint arXiv:1702.01992*, 2017.
- [Chen *et al.*, 2016] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *NeurIPS*, 29, 2016.
- [Chen *et al.*, 2020] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607. PMLR, 2020.
- [Chen *et al.*, 2021] Yanbei Chen, Yongqin Xian, A Koepke, Ying Shan, and Zeynep Akata. Distilling audio-visual knowledge by compositional contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7016–7025, 2021.
- [Chen *et al.*, 2023] Mengxi Chen, Linyu Xing, Yu Wang, and Ya Zhang. Enhanced multimodal representation learning with cross-modal kd. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11766–11775, 2023.
- [Daunhawer *et al.*, 2023] Imant Daunhawer, Alice Bizeul, Emanuele Palumbo, Alexander Marx, and Julia E Vogt. Identifiability results for multimodal contrastive learning. *arXiv preprint arXiv:2303.09166*, 2023.
- [Demšar, 2006] Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine learning research*, 7:1–30, 2006.
- [Feichtenhofer *et al.*, 2022] Christoph Feichtenhofer, Yang-hao Li, Kaiming He, et al. Masked autoencoders as spatiotemporal learners. *NeurIPS*, 35:35946–35958, 2022.
- [Huang *et al.*, 2021] Yu Huang, Chenzhuang Du, Zihui Xue, Xuanyao Chen, Hang Zhao, and Longbo Huang. What makes multi-modal learning better than single (provably). *NeurIPS*, 34:10944–10956, 2021.
- [Kiela *et al.*, 2020] Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ringshia, and Davide Testuggine. The hateful memes challenge: Detecting hate speech in multimodal memes. *NeurIPS*, 33:2611–2624, 2020.
- [Kim *et al.*, 2021] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision. In *ICML*, pages 5583–5594. PMLR, 2021.
- [Kingma and Ba, 2015] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, San Diego, CA, USA, 2015.
- [Kingma and Welling, 2014] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In Yoshua Bengio and Yann LeCun, editors, *ICLR*, 2014.
- [Lee and Van der Schaar, 2021] Changhee Lee and Mihaela Van der Schaar. A variational information bottleneck approach to multi-omics data integration. In *International Conference on Artificial Intelligence and Statistics*, pages 1513–1521. PMLR, 2021.
- [Lee *et al.*, 2023a] Seungyeol Lee, Taeho Kim, and Jae-Pil Heo. Cross-loss pseudo labeling for semi-supervised segmentation. *IEEE Access*, 2023.
- [Lee *et al.*, 2023b] Yi-Lun Lee, Yi-Hsuan Tsai, Wei-Chen Chiu, and Chen-Yu Lee. Multimodal prompting with missing modalities for visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14943–14952, 2023.
- [Li *et al.*, 2021] Shuang Li, Wenxuan Ma, Jinming Zhang, Chi Harold Liu, Jian Liang, and Guoren Wang. Meta-reweighted regularization for unsupervised domain adaptation. *IEEE Transactions on Knowledge and Data Engineering*, 35(3):2781–2795, 2021.
- [Liang *et al.*, 2023] Paul Pu Liang, Chun Kai Ling, Yun Cheng, Alex Obolenskiy, Yudong Liu, Rohan Pandey, Alex Wilf, Louis-Philippe Morency, and Ruslan Salakhutdinov. Multimodal learning without labeled multimodal data: Guarantees and applications. *arXiv preprint arXiv:2306.04539*, 2023.
- [Ma *et al.*, 2021] Mengmeng Ma, Jian Ren, Long Zhao, Sergey Tulyakov, Cathy Wu, and Xi Peng. Smil: Multimodal learning with severely missing modality. In *Proceedings of the AAAI Conference on Artificial Intelligence*, number 3 in 35, pages 2302–2310, 2021.
- [Ma *et al.*, 2022] Mengmeng Ma, Jian Ren, Long Zhao, Davide Testuggine, and Xi Peng. Are multimodal transformers robust to missing modality? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18177–18186, 2022.
- [Peng *et al.*, 2005] Hanchuan Peng, Fuhui Long, and Chris Ding. Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy. *IEEE Transactions on pattern analysis and machine intelligence*, 27(8):1226–1238, 2005.

- [Peng *et al.*, 2022] Xiaokang Peng, Yake Wei, Andong Deng, Dong Wang, and Di Hu. Balanced multimodal learning via on-the-fly gradient modulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8238–8247, 2022.
- [Poklukar *et al.*, 2022a] Petra Poklukar, Vladislav Polianskii, Anastasiia Varava, Florian T Pokorny, and Danica Kragic. Delaunay component analysis for evaluation of data representations. In *International Conference on Learning Representations; Apr 25th through Fri the 29th, 2022 (online)*, 2022.
- [Poklukar *et al.*, 2022b] Petra Poklukar, Miguel Vasco, Hang Yin, Francisco S Melo, Ana Paiva, and Danica Kragic. Geometric multimodal contrastive representation learning. In *ICML*, pages 17782–17800. PMLR, 2022.
- [Radford *et al.*, 2021] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763. PMLR, 2021.
- [Su *et al.*, 2021] Mingyue Su, Guanghua Gu, Xianlong Ren, Hao Fu, and Yao Zhao. Semi-supervised knowledge distillation for cross-modal hashing. *IEEE Transactions on Multimedia*, 2021.
- [Tsai *et al.*, 2018] Yao-Hung Hubert Tsai, Paul Pu Liang, Amir Zadeh, Louis-Philippe Morency, and Ruslan Salakhutdinov. Learning factorized multimodal representations. *arXiv preprint arXiv:1806.06176*, 2018.
- [Tsai *et al.*, 2019] Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences. In *Proceedings of the conference. Association for Computational Linguistics. Meeting*, volume 2019, page 6558. NIH Public Access, 2019.
- [Van der Maaten and Hinton, 2008] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [Wang and Isola, 2020] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th ICML*, volume 119 of *Proceedings of Machine Learning Research*, pages 9929–9939. PMLR, 13–18 Jul 2020.
- [Wang *et al.*, 2015] Xin Wang, Devinder Kumar, Nicolas Thome, Matthieu Cord, and Frederic Precioso. Recipe recognition with large multimodal food dataset. In *2015 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*, pages 1–6. IEEE, 2015.
- [Wang *et al.*, 2020] Qi Wang, Liang Zhan, Paul Thompson, and Jiayu Zhou. Multimodal learning with incomplete modalities by knowledge distillation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1828–1838, 2020.
- [Wang *et al.*, 2023] Yuanzhi Wang, Zhen Cui, and Yong Li. Distribution-consistent modal recovering for incomplete multimodal learning. In *CVPR*, pages 22025–22034, 2023.
- [Williams and Beer, 2010] Paul L Williams and Randall D Beer. Nonnegative decomposition of multivariate information. *arXiv preprint arXiv:1004.2515*, 2010.
- [Woo *et al.*, 2023] Sangmin Woo, Sumin Lee, Yeonju Park, Muhammad Adi Nugroho, and Changick Kim. Towards good practices for missing modality robust action recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, number 3 in 37, pages 2776–2784, 2023.
- [Wu and Goodman, 2018] Mike Wu and Noah Goodman. Multimodal generative models for scalable weakly-supervised learning. *NeurIPS*, 31, 2018.
- [Xu *et al.*, 2024] Jie Xu, Shuo Chen, Yazhou Ren, Xiaoshuang Shi, Hengtao Shen, Gang Niu, and Xiaofeng Zhu. Self-weighted contrastive learning among multiple views for mitigating representation degeneration. *NeurIPS*, 36, 2024.
- [Zadeh *et al.*, 2016] Amir Zadeh, Rowan Zellers, Eli Pincus, and Louis-Philippe Morency. Multimodal sentiment intensity analysis in videos: Facial gestures and verbal messages. *IEEE Intelligent Systems*, 31(6):82–88, 2016.
- [Zadeh *et al.*, 2018a] Amir Zadeh, Paul Pu Liang, Soujanya Poria, Prateek Vij, Erik Cambria, and Louis-Philippe Morency. Multi-attention recurrent network for human communication comprehension. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [Zadeh *et al.*, 2018b] AmirAli Bagher Zadeh, Paul Pu Liang, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. Multimodal language analysis in the wild: Cmu-mosei dataset and interpretable dynamic fusion graph. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2236–2246, 2018.
- [Zhang *et al.*, 2022a] Shu Zhang, Ran Xu, Caiming Xiong, and Chetan Ramaiah. Use all the labels: A hierarchical multi-label contrastive learning framework. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16660–16669, 2022.
- [Zhang *et al.*, 2022b] Xinyu Zhang, Zhiwei Li, Yan Gong, Dafeng Jin, Jun Li, Li Wang, Yanzhang Zhu, and Huaping Liu. Openmpd: An open multimodal perception dataset for autonomous driving. *IEEE Transactions on Vehicular Technology*, 71(3):2437–2447, 2022.
- [Zhang *et al.*, 2023] Shuo Zhang, Jiaoqiao Zhang, Biao Tian, Thomas Lukasiewicz, and Zhenghua Xu. Multi-modal contrastive mutual learning and pseudo-label re-learning for semi-supervised medical image segmentation. *Medical Image Analysis*, 83:102656, 2023.

A Related Works

Semi-supervised Multimodal Learning: Several works acknowledge the challenge of fully labeled datasets in the multimodal literature and provide targeted solutions for specific applications [Aberdam *et al.*, 2022; Zhang *et al.*, 2023; Liang *et al.*, 2023]. For instance, in [Aberdam *et al.*, 2022], the authors tackle the semi-supervised scenario in scene text recognition by enforcing consistency between weakly augmented pseudo-labels and strongly augmented views. Authors in [Zhang *et al.*, 2023] propose an area-similarity contrastive loss for medical image segmentation, leveraging cross-modal information to enhance representations of unlabeled data. Liang et al. [Liang *et al.*, 2023] derive two lower bounds of multimodal interaction from an information-theoretic perspective, applicable for pre-analysis of multimodal interaction effects. A common technique in general semi-supervised learning is loss reweighting based on pseudolabel uncertainty, similar to our PU loss. These methods aim to mitigate confirmation bias and are widely used in unsupervised domain adaptation [Li *et al.*, 2021], particularly in image processing applications [Lee *et al.*, 2023a]. However, these efforts primarily focus on semi-supervised scenarios in specific tasks and certain modalities (e.g., text-images), limiting their applicability to general cases.

Missing modalities: Many multimodal fusion methods rely on a complete set of modalities, but deployment settings often lack such ideal conditions, leading to adverse effects when using these strategies [Wang *et al.*, 2020; Ma *et al.*, 2022]. To address this challenge, some approaches aim to create models resilient to missing modalities [Ma *et al.*, 2021; Ma *et al.*, 2022; Poklukar *et al.*, 2022b; Woo *et al.*, 2023; Lee *et al.*, 2023b]. For instance, Wang et al. [Wang *et al.*, 2020] optimize training by considering incomplete data samples to generate unimodal teachers guiding a multimodal student. Smil [Ma *et al.*, 2021] approximates latent features of modality-incomplete data using Bayesian meta-learning. GMC [Poklukar *et al.*, 2022b] preserves geometric alignment in multimodal representations, enabling unimodal representations to substitute for absent representations of other modalities. ActionMAE, inspired by the masked autoencoder idea [Feichtenhofer *et al.*, 2022], learns to predict the latent representation of a missing modality by randomly dropping its feature token and learning to reconstruct it. Despite success in certain scenarios, these frameworks often rely on labeled signals, implicitly or explicitly, in the training dataset, limiting their general applicability.

B Robult supplementary details

B.1 Minimizing Conditional Entropy with Latent Reconstruction Error

As explained in the primary text, our approach to achieve Objective 2.2 involves a reconstruction procedure with two components: the reconstruction module $r^i(U^i, Z^i) = \tilde{H}^i$ and the latent reconstruction loss \mathcal{L}_{rec} . This procedure is illustrated in Figure 6. It is important to note that these reconstruction modules $r^i(\cdot)$ are used exclusively during the learning process to optimize individual branches $g^i(\cdot)$, incurring no additional overhead during the evaluation or deployment stages. As this reconstruction is carried out in the latent space, the module $r^i(\cdot)$ can be uniformly designed, irrespective of the characteristics of input modalities. In our Robult design, $r^i(\cdot)$ is simply a two-layered MLP with ReLU activation in the middle, applied across all five datasets.

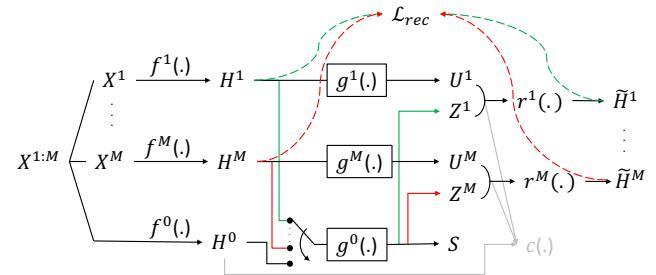


Figure 6. Reconstruction procedure of Robult. This procedure only applied in training stage.

In this section, we would derive the lower bound of mutual information between fused latent S and unimodal representation Z^i as a Positive-Unlabelled learning objective, which relax the assumption about full presence of labels in training dataset. This derivation explains Result 2 in the main manuscript.

B.2 Maximizing Mutual Information with Soft Positive-Unlabeled Contrastive Learning

The ultimate goal is to maximize the following MI quantity:

$$\mathcal{I}(S, Z^i) = d_{KL}(p_{S, Z^i} || p_S \otimes p_{Z^i})$$

This essentially means that the KL divergence between the joint distribution p_{S, Z^i} and the product of marginal distribution $p_S \otimes p_{Z^i}$ should be maximized. As defined in the main manuscript, F is the flag indicator denotes whether a couplet (s, z^i) is sampled from the joint distribution p_{S, Z^i} ($F = 1$) or from product of marginal distribution $p_S \otimes p_{Z^i}$ ($F = 0$):

$$p(S, Z^i | F = 1) = p_{S, Z^i}; \quad p(S, Z^i | F = 0) = p_S \otimes p_{Z^i}; \quad (11)$$

Applying Bayes' rule, the posterior for $F = 1$ is given by:

$$\begin{aligned} p(F = 1 | S, Z^i) &= \frac{p(S, Z^i | F = 1)p(F = 1)}{p(S, Z^i)} \\ &= \frac{p(S, Z^i | F = 1)p(F = 1)}{p(S, Z^i | F = 1)p(F = 1) + p(S, Z^i | F = 0)p(F = 0)} \\ &= \frac{p_{S, Z^i} \cdot p(F = 1)}{p_{S, Z^i} \cdot p(F = 1) + p_S \otimes p_{Z^i} \cdot p(F = 0)}. \end{aligned} \quad (12)$$

Applying logarithm operation on both side of Equation

$$\begin{aligned} \log p(F = 1|S, Z^i) &= -\log \left(1 + k \frac{p_S \otimes p_{Z^i}}{p_{S, Z^i}} \right) \\ &\leq -\log k + \log \frac{p_{S, Z^i}}{p_S \otimes p_{Z^i}}, \end{aligned} \quad (13)$$

in which

$$k = \frac{p(F = 0)}{p(F = 1)}. \quad (14)$$

Taking expectation w.r.t p_{S, Z^i} (or $p(S, Z^i|F = 1)$), we can bound the mutual information as

$$\mathcal{I}(S, Z^i) \geq \mathbb{E}_{p(S, Z^i|F=1)} \log p(F = 1|S, Z^i) + \log k \quad (15)$$

Here, the true distribution $p(F = 1|S, Z^i)$ is unknown, so we approximate it with a well-established non-parametric model $v : S \times Z^i \rightarrow [0, 1]$ [Chen *et al.*, 2021; Chen *et al.*, 2023]:

$$\mathcal{I}(S, Z^i) \geq \mathbb{E}_{p(S, Z^i|F=1)} \log v(S, Z^i) + \log k$$

where

$$\begin{aligned} v(s_j, z_k^i) &= \frac{\phi(s_j, z_k^i)}{\sum_h^B \phi(s_j, z_h^i)}; \\ \phi(s_j, z_k^i) &= \exp(\langle s_j; z_k^i \rangle / \tau). \end{aligned} \quad (16)$$

In addition, let c be the number of underlying classes and assume the labels are uniformly distributed, we have the probability that a couplet is sharing a label is $p(f_k = 1) = \frac{1}{c^2}$. Within mini-batch of size B , consider the scenario in which the number of positive couplets B_p is greater than the number of negative ones B_n (hence $B_p > \frac{\binom{B}{2}}{2} = \frac{\tilde{B}}{2}$):

$$\begin{aligned} p(B_p) &= \binom{\tilde{B}}{B_p} \cdot \frac{1}{c^{2B_p}} \\ &\leq \binom{\tilde{B}}{\frac{\tilde{B}}{2}} \cdot \frac{1}{c^{\tilde{B}}} \end{aligned}$$

It should be noted that this possibility $p(B_p)$ is upper-bounded by a small quantity given $B > 1$ and $c \geq 2$ (smaller than 0.1 in case $B = 8$ and $c = 2$), and get smaller when B and c increase. Intuitively, the possibility that the positive couplets outnumber the negative ones is negligible, hence, it normally hold true that:

$$\log k = \log \frac{p(F = 0)}{p(F = 1)} \geq 0.$$

With this realization, Result 2 can be derived from Eq. 16 as follow:

$$\begin{aligned} \mathcal{I}(S, Z^i) &\geq \mathbb{E}_{p(S, Z^i|F=1)} \log v(S, Z^i) + \log k \\ &\geq \mathbb{E}_{p(S, Z^i|F=1)} \log v(S, Z^i) \\ &= \mathbb{E}_{p_{S, Z^i}} \log v(S, Z^i). \end{aligned} \quad (17)$$

B.3 Training Strategy

We employ an end-to-end training pipeline that can process two objectives 2.2 and 2.1 independently, as demonstrated by Algorithm 1. In general, we selectively perform gradient calculations on different modules of Robult based on the specific losses. This selection process offers dual benefits: (1) minimal processing overhead with a single forward pass, and (2) effective restriction of the losses' impact only on their target modules.

Algorithm 1 Robult training strategy

Input:

- ▷ Training dataset \mathcal{D}_{train}
- ▷ Robult framework \mathcal{RB}
- ▷ Optimizer \mathcal{O}

function *ParametersToggle*(f : flag-variable)

- if** $f = 0$ **then**
- Toggle all \mathcal{RB} parameters to require gradient calculation
- else if** $f = 1$ **then**
- Toggle all \mathcal{RB} parameters to NOT require gradient calculation
- Toggle all $g^i(\cdot)$ parameters ($i = 1, \dots, M$) to require gradient calculation
- else if** $f = 2$ **then**
- Toggle all \mathcal{RB} parameters to NOT require gradient calculation
- Toggle all $f^i(\cdot), g^i(\cdot)$ parameters ($i = \overline{0, M}$) to require gradient calculation
- end if**

end function

for B_i, Y_i **in** \mathcal{D}_{train} **do**

- ▷ single forward pass
- $\tilde{Y}_i, H_i, Z_i, U_i, S = \mathcal{RB}(B_i)$
- ▷ loss calculations
- $l_{cls} = \mathcal{L}_{cls}(\tilde{Y}_i, Y_i)$
- $l_{rec} = \mathcal{L}_{rec}(H_i, Z_i, U_i)$
- $l_{(u)lb} = \mathcal{L}_{(u)lb}(Z_i, S)$
- ▷ gradient calculations
- Call *ParametersToggle*($f = 1$); backward with l_{rec}
- Call *ParametersToggle*($f = 2$); backward with l_{lb} and l_{ulb}
- Call *ParametersToggle*($f = 0$); backward with l_{cls}
- ▷ single backward pass
- Optimizer \mathcal{O} update \mathcal{RB} parameters with above gradient infomation

end for

B.4 Robult's Complexity Analysis

Robult framework is built on two main types of modules: individual branch modules - $g^i(\cdot)$ ($i = 0, M$), and reconstruction modules - $r^i(\cdot)$ ($i = 1, \dots, M$). For the projectors and the fusion module, denoted as $f^i(\cdot)$ ($i = \overline{0, M}$), we adopt designs from previous studies. In this section, we analyze the complexity of our two proposed modules, which operate in latent spaces and have straightforward designs.

Individual branches

Since all $g^i(\cdot)$'s are working with the same input latent space (all raw modalities are projected to the same space), we unify the design of $g^i(\cdot)$'s to be identical across different modalities. Specifically, $g^i(\cdot)$ are constituted by multiple Fully Connected layers, with middle *ReLU* activations; the last layer of $g^i(\cdot)$ involve no activation, but a L2 normalization operation. Below are the table of hyper-parameters involved in the analysis.

Notation	Description
M	number of modalities
L	number of FC layers
d_i	hidden dimension of i^{th} layer's output
d_0	input dimension

Table 4. $g^i(\cdot)$ related hyper-parameters

Time Complexity. Assume a single operation can be performed in unit time ($\mathcal{O}(1)$). We have the calculation for number of operations in a forward pass as follows.

Within the i^{th} FC layer:

$$d_{i-1} * d_i + d_i,$$

Over L layers:

$$\sum_{i=1}^L d_{i-1} * d_i + d_i.$$

In our implementations, we choose the same dimensions for all hidden outputs (same $d = d_i \forall i = \overline{1, L}$), and there are $M+1$ modules $g^i(\cdot)$. With this, the total number of operation is:

$$(M+1) \sum_{i=1}^L d_{i-1} * d_i + d_i = (M+1) * L * d * (d+1) = \mathcal{O}(M * L * d^2)$$

By utilizing matrix product and GPU acceleration, d^2 operations can in fact be performed in $\mathcal{O}(1)$ time, make the whole time complexity for individual branches be $\mathcal{O}(M * L)$, which is linearly scaled with M .

Space Complexity. Regarding the space complexity, within i^{th} layer, beside the need for storing parameter matrix of size $(d_{i-1} + 1) \times d_i$, output after performing *ReLU* activation are also stored to later perform back-propagation. Hence, the total number of stored parameters is:

$$(d_{i-1} + 1) * d_i + d_i = (d_{i-1} + 2) * d_i.$$

Following similar derivation with L layers and $M + 1$ branches, replacing $d = d_i \forall i = \overline{1, L}$, we have the total space complexity is:

$$(M + 1) * L * (d + 2) * d = \mathcal{O}(M * L * d^2).$$

Reconstruction modules

For these reconstruction modules $r^i(\cdot)$'s, we also adopt a similar design patterns as that of individual branches $g^i(\cdot)$. The

only differences are the dimension of input for first FC layer ($2d$), which corresponding to the concatenation of $g^0(\cdot)$'s and $g^i(\cdot)$'s outputs.

Time complexity. With that intuition, as we have M branches and L layers, the total number of calculations is:

$$\begin{aligned} & M * [(2d + 1) * d + (L - 1) * (d + 1) * d] \\ & = M * [d^2 + L * (d + 1) * d] = \mathcal{O}(M * L * d^2) \end{aligned}$$

Reducing d^2 operations to $\mathcal{O}(1)$ time complexity, the same result as observed with $g^i(\cdot)$'s are observed - $\mathcal{O}(M * L)$.

Space complexity. The total number of stored parameters is:

$$\begin{aligned} & M * [(2d + 2) * d + (L - 1) * (d + 2) * d] \\ & = M * [d^2 + L * (d + 2) * d] = \mathcal{O}(M * L * d^2). \end{aligned}$$

In conclusion, all the proposed modules of Robult are linearly scaled (both in time and space), with the number of modalities M .

Computational Time Quantitative Result

Robult	GMC	ActionMAE
<i>CMU-MOSI:</i>		
1.08 GFLOPS	1.05 GFLOPS	1.13 GFLOPS
507.44 MMACs	492 MMACs	526.36 MMACs
1.46 M	1.16 M	11.55 M
<i>CMU-MOSEI:</i>		
1.09 GFLOPS	1.06 GFLOPS	1.15 GFLOPS
508.7 MMACs	493.26 MMACs	526.36 MMACs
1.41 M	1.17 M	11.56 M
<i>Hateful Memes:</i>		
32.81 MFLOPs	26.56 MFLOPS	61.37 MFLOPS
16.39 MMACs	13.27 MMACs	30.48 MMACs
888.32 K	674.05 K	1.04 M

Table 5. Computational times of different methods on different datasets.

To further substantiate our results, we measured FLOPs and MACs for several datasets we utilized, comparing them with our current baselines (Table 5). For a fair comparison, we kept all unimodal projectors the same. The results suggest that our method introduces slight overheads compared to GMC, but remains faster than ActionMAE, while significantly outperforming both methods in downstream performance.

C Implementation details

C.1 Environment Settings

All implementations and experiments are conducted on a single machine equipped with the following hardware configuration: a 6-core Intel Xeon CPU paired with 2 NVIDIA A100 GPUs for accelerated training.

Our codebase predominantly utilizes the *PyTorch 2.0* framework, including the *Pytorch-AutoGrad*, for deep learning model design and calculations. Additionally, we leverage utilities from *Scikit-learn*, *Pandas*, and *Matplotlib* to support various functionalities in our experiments. The original codebase for Robult will be made publicly available upon publication.

C.2 Reproduction and Adaptation

In both Robult and the Unimodal baselines, we modify the architecture of the projectors $f^i(\cdot)$ ($i = \overline{0, M}$) while keeping $g^i(\cdot)$ ($i = \overline{0, M}$) as simple as possible. For all testing datasets, the unimodal branches $g^i(\cdot)$ consist of a simple Fully Connected layer followed by L_2 normalization. In contrast, $g^0(\cdot)$ has a higher representation capacity with two Fully Connected layers and ReLU activation. In the case of GMC [Poklukar *et al.*, 2022b] and ActionMAE [Woo *et al.*, 2023], the same architecture of the projectors is adopted as Robult to ensure a fair comparison, with the remaining designs being directly inherited from the original codebases. For Prompt-Trans [Lee *et al.*, 2023b], we keep all the architecture designs intact and only change the datasets' settings to semi-supervised and missing modalities scenarios. Additional information about the baselines should be best referenced from their original works.

CMU-MOSI and CMU-MOSEI datasets. We follow the settings in [Poklukar *et al.*, 2022b], which involve temporally-aligned versions of these datasets generated with [Zadeh *et al.*, 2018a], with additional adaptions for semi-supervised and missing modalities scenarios. Specifically, multimodal Transformer [Tsai *et al.*, 2019] is adopted as the joint-modality encoder $f^0(\cdot)$ for our model and all state-of-the-art baselines; single-layer GRUs are adopted as unimodal projectors $f^i(\cdot)$. The latent space dimension is set as 60, and all methods are trained in 40 epochs with Adam optimizer [Kingma and Ba, 2015] at the learning rate of 10^{-3} .

MM-IMDb, UPMC Food-101 and Hateful Memes datasets. For the classification tasks associated with the MM-IMDb, UPMC Food-101, and Hateful Memes datasets, we initially generate text and visual embeddings offline using a pretrained ViLT framework [Kim *et al.*, 2021]. Subsequently, all models are trained using these embeddings instead of raw data. This specific procedure is intentionally conducted to ensure fair evaluation, as Prompt-Trans [Lee *et al.*, 2023b] functioning also involves the same frozen ViLT framework in their training process. We also follow Prompt-Trans [Lee *et al.*, 2023b] for the preprocessing procedures of raw texts and images. For all methods, the offline embedding space's dimension is fixed at 784 and further condensed into a 128-dimensional hidden latent space. Additionally, with this setting, we choose the projectors $f^i(\cdot)$ ($i = \overline{0, M}$) as simple Fully Connected layers.

D Additional empirical results and analysis

In this section, we present additional empirical results and analysis to study the behavior of Robult and baselines in dif-

ferent extended settings.

D.1 Extended modalities missing scenarios

In Table 6, we present the comprehensive performance of different frameworks when provided access to all combinations of input modalities on CMU-MOSI and CMU-MOSEI datasets. This table extends the information presented in Table 1 in the main text. In this experiment, to report the performances of Unimodal baselines and Robult when provided with two modalities, we simply take the mean of the outputs generated by providing these frameworks with single modalities. It is important to note that we do not draw conclusions on the best strategy for merging unimodal results. Despite this, using this simple strategy, Robult consistently produces the best results in most scenarios, highlighting its performance consistency across different missing modality scenarios.

D.2 Extended semi-supervised scenarios

This experiment is designed to observe the behavior of models when exposed to varying amounts of labeled information. In addition to the 5% labeled ratio setting covered in the main text, we additionally evaluate all methods with 50% labeled ratio and ideal supervised settings:

- **Semi-supervised learning with 50% labelled data.** Table 7 summarizes the results of the 50% labeled setting with CMU-MOSI and Hateful Memes datasets. As indicated, all methods effectively leverage the increased label signal, resulting in improved performance. However, it's noteworthy that Robult demonstrates its superiority by outperforming other methods on both datasets across various metrics.
- **Supervised learning.** With this scenario, we evaluate all methods in an ideal case where fully labelled training dataset is available. Similar to the previous setting, all method further enhance their performance given more labeled data. Robult suggest the consistency by outperforming other baselines in most recorded metrics.

To clearly illustrate the performance improvement of Robult and the baselines in each scenario, we provide visualizations of Pearson correlation for CMU-MOSI and AUROC for Hateful Memes in Figure 7. Among all methods, Robult demonstrates the best stability and consistency in performance, regardless of input modalities or label ratios.

D.3 Extended comparison with recent frameworks

Baselines. We adopt two recent approaches utilizing Contrastive Loss [Xu *et al.*, 2024] and reconstruction strategy [Wang *et al.*, 2023] for more comprehensive comparison of Robult with existing State-of-the-art frameworks. Their original codebases are slightly adjusted for semi-supervised settings, and the dimensions of the latent space are aligned with Robult's (60) to minimize bias in the comparison.

Settings and Result. We evaluate the models' performance using a 5% semi-supervised task with the CMU-MOSI and

Metrics	CMU-MOSI				CMU-MOSEI			
	Unimodal	GMC	ActionMAE	Robult	Unimodal	GMC	ActionMAE	Robult
<i>Text Modality:</i>								
MAE	1.41	1.407	1.476	1.397	0.81	0.815	1.115	0.784
Corr	0.137	0.14	0.066	0.144	0.383	0.346	0.136	0.459
F1	0.551	0.559	0.535	0.578	0.717	0.716	0.614	0.739
Acc	0.553	0.562	0.47	0.569	0.712	0.708	0.603	0.732
<i>Audio Modality:</i>								
MAE	1.576	1.518	1.546	1.415	0.842	0.836	1.215	0.825
Corr	0.041	-0.065	0.046	0.085	0.111	0.193	0.101	0.221
F1	0.512	0.457	0.508	0.539	0.618	0.642	0.634	0.679
Acc	0.496	0.46	0.467	0.535	0.599	0.63	0.543	0.65
<i>Vision Modality:</i>								
MAE	1.451	1.497	1.511	1.425	0.891	0.839	1.127	0.826
Corr	0.044	-0.07	-0.03	0.086	0.163	0.2	0.104	0.201
F1	0.585	0.446	0.511	0.593	0.637	0.621	0.594	0.647
Acc	0.425	0.449	0.514	0.522	0.624	0.62	0.561	0.632
<i>Text+Audio Modalities:</i>								
MAE	1.485	1.442	1.521	1.401	0.765	0.813	1.007	0.762
Corr	0.131	0.05	0.089	0.141	0.418	0.352	0.202	0.439
F1	0.528	0.491	0.507	0.563	0.729	0.728	0.624	0.733
Acc	0.512	0.493	0.508	0.546	0.713	0.671	0.62	0.717
<i>Text+Vision Modalities:</i>								
MAE	1.486	1.465	1.489	1.415	0.861	0.822	1.003	0.788
Corr	0.144	0.044	0.086	0.146	0.325	0.325	0.198	0.399
F1	0.514	0.487	0.501	0.58	0.718	0.722	0.623	0.718
Acc	0.492	0.491	0.506	0.534	0.688	0.687	0.619	0.704
<i>Audio+Vision Modalities:</i>								
MAE	1.432	1.499	1.534	1.426	0.824	0.824	1.173	0.812
Corr	0.014	-0.075	-0.035	0.091	0.214	0.233	0.147	0.244
F1	0.492	0.445	0.55	0.581	0.748	0.663	0.637	0.663
Acc	0.486	0.448	0.453	0.527	0.638	0.633	0.623	0.64
<i>Full Modalities:</i>								
MAE	1.394	1.47	1.496	1.392	0.783	0.819	1.103	0.779
Corr	0.186	0.101	-0.092	0.247	0.364	0.328	0.337	0.504
F1	0.597	0.497	0.553	0.657	0.73	0.693	0.694	0.744
Acc	0.594	0.498	0.477	0.63	0.729	0.688	0.643	0.741

Table 6. Full performance of different frameworks on CMU-MOSI and CMU-MOSEI Dataset.

CMU-MOSEI datasets, testing all possible combinations of modalities input. Table 9 summarizes the results of this study. As shown, Robult consistently outperforms the two frameworks in most scenarios. This experiment further highlights Robult’s robustness in semi-supervised settings and when modalities are missing.

D.4 Extended ablation studies on Robult design

Setting. In this analysis, our goal is to understand the contributions of our applied strategies to overall Robult’s performance. Specifically, we adopt several ablation studies:

- **Removal of Unimodal branches** $g^i(\cdot)(i = 1 \dots M)$: The output of the shared branch $g^0(\cdot)$ is directly fed into the classifier to yield the final result. The remaining framework is trained normally with the soft PU loss and downstream task loss.

- **Soft-PU Loss Ablation - Removal of weighting scheme:** Uniform weight is adopted instead of our proposed dynamic weighting scheme.

- **Soft-PU Loss Ablation - Removal of pseudo labeling:** All unlabeled samples are considered negatives, resemble normal constrastive learning scheme.

Result. The results, presented in Table 10, indicate an overall performance decrease across all modalities on two tested datasets. Specifically, with removal of unimodal branches, in the case of a small dataset with few labeled samples (CMU-MOSI), this ablation causes some weaker modalities to fail in generating beneficial representations during learning. Similar patterns are captures with ablations of Soft P-U loss. The results indicate a consistent decrease in performance across both variations and two test datasets. This analysis empirically supports the effectiveness of our soft PU loss.

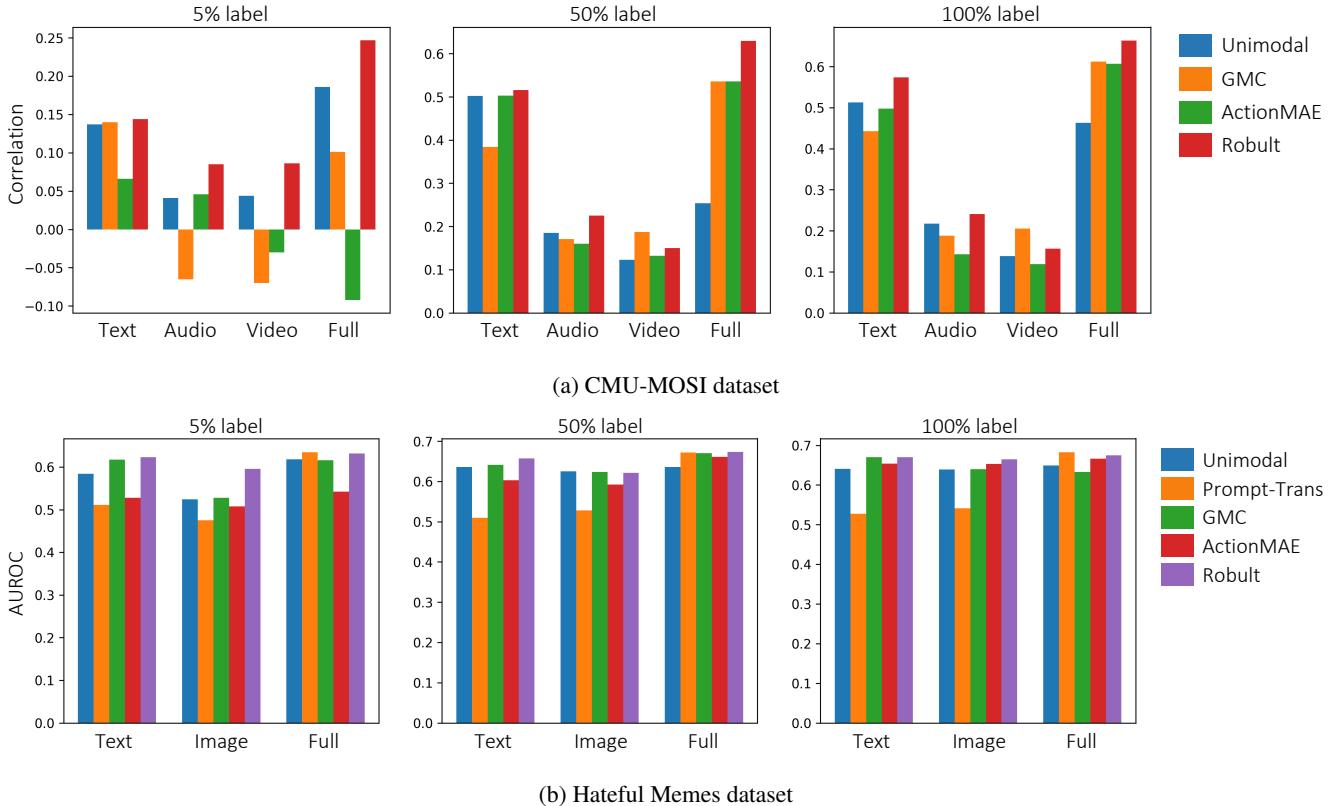


Figure 7. Models’ performances when being exposed to different label ratios.

D.5 Extended ablation study regarding choice of RBF Kernel

In this analysis, our goal is to understand the role of our weighting scheme in the Soft P-U Loss. We compare two distinct weighting mechanisms to evaluate how closely a positive candidate matches the true positive pair, and then contrast these mechanisms against our initial choice of the RBF Kernel.

Setting. The two new weighting mechanisms are designed based on normalized distances, with the difference lying in the choice of distance measures $\delta(\cdot)$ (here Euclidean and Manhattan distances). Specifically, within a mini-batch B , given the reference proximity ϕ_{ref} and the proximity ϕ_i of the positive candidate that need to be weighted, we calculated the weight as follow:

$$w_i = 1 - \tilde{d}_i;$$

$$\tilde{d}_i = \frac{\delta(\phi_i, \phi_{ref})}{\max_B \delta(\phi_j, \phi_{ref})}.$$

Result. We refer to the variant using a Manhattan distance-based strategy as *Robult - L1*, and the one utilizing Euclidean measures as *Robult - L2*. These two variants are evaluated against the original RBF-based model in a 5% semi-supervised task with the CMU-MOSI and CMU-MOSEI datasets. The comprehensive results are presented in Table 11. Generally, we observe minor differences in performance

among the weighting schemes. While the RBF approach yields the most consistent results across various input combinations for these datasets, we do not declare it the definitive best weighting method. We believe further research is needed to identify the most appropriate strategy for the dataset of interest.

D.6 Alignment and Uniformity Analysis

We assess the learned representations Z^i and S after the training process with soft P-U loss, via two qualities - Alignment and Uniformity [Wang and Isola, 2020].

Figure 8 provides a comprehensive analysis of the learned representations generated by Robult using both unimodal and multimodal inputs on the Hateful Memes testing set. On the left, the Frobenius-norm distance histograms of positive pairs within the test dataset indicate that the representations generated with all the modalities have the smallest mean distances, and as the distances increase, their corresponding density decreases. While not as compact as the representations with full modalities input, positive pairs’ representations generated with unimodal input still exhibit low mean distances and good histogram shapes. Furthermore, to analyze the uniformity characteristics of the learned representations, we follow the process outlined in [Wang and Isola, 2020] and show the result on the right of Figure 8. The learned representations are projected into \mathbb{R}^2 using t-SNE [Van der Maaten and Hinton, 2008], and the output feature distributions are visualized us-

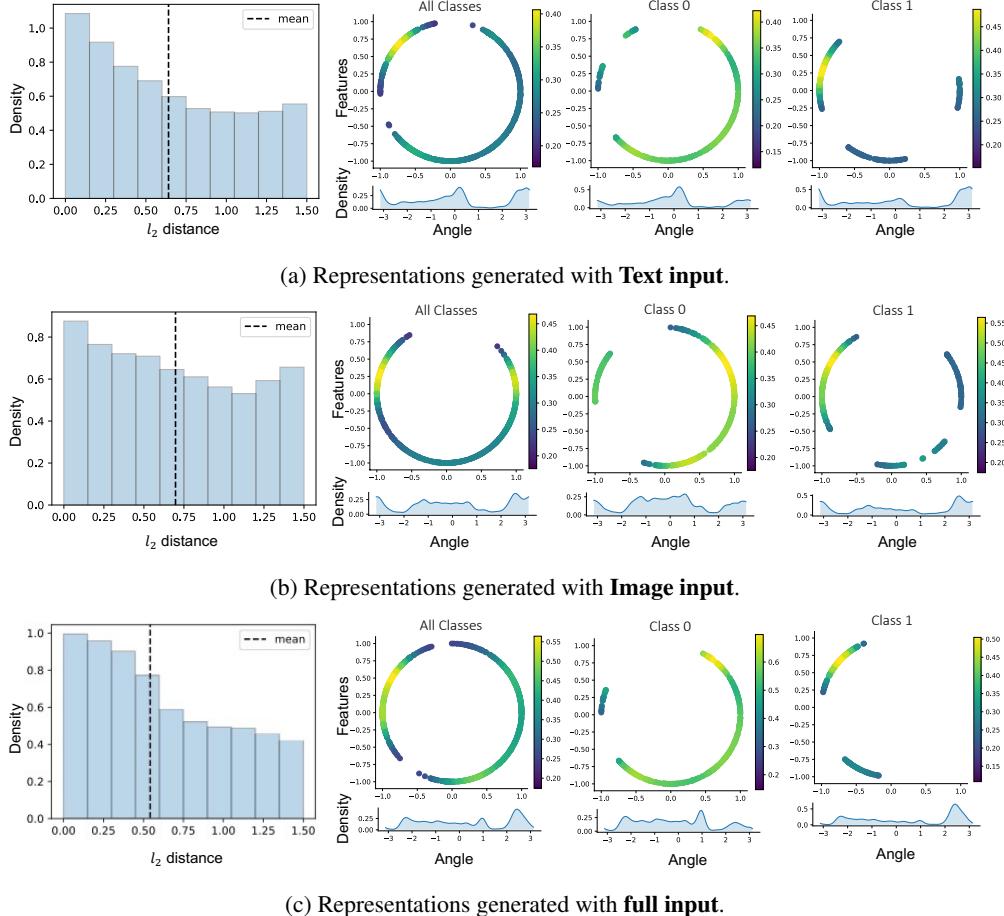


Figure 8. Alignment and Uniformity analysis on representations of Hateful Memes test dataset, generated by Robult.

ing Gaussian kernel density estimation (KDE) along with von Mises-Fisher (vMF) KDE for angles ($\arctan2(y; x)$). As suggested by these figures, Robult’s representations demonstrate uniform characteristics on the entire test set as well as good clustering between classes. Specifically, representations of different classes reside on different segments of the unit circle and form separated clusters in Figure 10. The level of separation for different classes with different input modalities correlates well with the actual quantitative results, as shown in Table 2. Additional clustering comparison between different method can be found in Appendix D.9.

D.7 Mutual Information Maximization Analysis

As stated in our main text, the necessity to model the objective of learning unimodal representations to maximize mutual information with a lower bound arises because mutual information cannot be precisely calculated. This is due to the changing values of the variables over time and the discrete nature of the datasets. To verify the effectiveness of our proposed method, we adopt the histogram-based method in [Peng *et al.*, 2005] to approximate MI between two variables after the training process with and without our soft PU loss (Table 12). The result suggest two important points:

- With our soft PU loss, the mutual information of all unimodal representations with the fused representation increase significantly.
- The entropy of the fused representation also increases with the use of our loss, suggesting that the fused representation also get enriched after training with the soft PU loss.

D.8 Soft label quality Analysis

We acknowledge that the quality of pseudo-labels is crucial for effective model training. This is why we incorporate our weighting scheme into the Positive-Unlabeled (PU) contrastive loss, considering the stochastic and unstable nature of pseudo-labels. This approach helps to reduce the impact of noisy pseudo-labels on the training process.

To demonstrate the effect of both pseudo-labels and our weighting strategy, we visualize the confusion matrix of pseudo-labels with and without the weighting scheme, compared to ground truth labels (Figure 9). This figure is plotted at epoch 20 of our training process using the CMU-MOSI dataset. The confusion matrix shows a strong correlation between pseudo-labels and ground truth labels, and the weighting scheme (removing all samples with weights below the

Modality	Metrics	Framework				
		Unimodal	Prompt-Trans	GMC	ActionMAE	Robult
Text	MAE	1.157	-	1.286	1.161	1.137
	Corr	0.502	-	0.384	0.503	0.516
	F1	0.712	-	0.642	0.701	0.721
	Acc	0.712	-	0.644	0.709	0.722
Audio	MAE	1.443	-	1.44	1.603	1.363
	Corr	0.185	-	0.171	0.16	0.225
	F1	0.563	-	0.556	0.566	0.569
	Acc	0.559	-	0.548	0.521	0.571
Video	MAE	1.514	-	1.458	1.406	1.429
	Corr	0.123	-	0.187	0.132	0.15
	F1	0.541	-	0.566	0.551	0.57
	Acc	0.49	-	0.56	0.535	0.561
Full	MAE	1.508	-	1.148	1.096	1.092
	Corr	0.254	-	0.536	0.536	0.63
	F1	0.566	-	0.709	0.728	0.761
	Acc	0.545	-	0.71	0.705	0.762
Text	AUROC	0.636	0.51	0.641	0.603	0.657
	Accuracy	0.583	0.508	0.592	0.556	0.589
Image	AUROC	0.625	0.528	0.624	0.592	0.621
	Accuracy	0.566	0.526	0.568	0.53	0.568
Full	AUROC	0.636	0.672	0.67	0.661	0.673
	Accuracy	0.57	0.594	0.596	0.573	0.604

Table 7. Semi-supervised learning with 50% labelled data on CMU-MOSI and Hateful Memes datasets.

Modality	Metrics	Framework				
		Unimodal	Prompt-Trans	GMC	ActionMAE	Robult
Text	MAE	1.126	-	1.233	1.108	1.066
	Corr	0.513	-	0.443	0.498	0.574
	F1	0.716	-	0.665	0.739	0.756
	Acc	0.717	-	0.667	0.749	0.753
Audio	MAE	1.421	-	1.414	1.569	1.392
	Corr	0.217	-	0.188	0.143	0.241
	F1	0.574	-	0.571	0.569	0.574
	Acc	0.553	-	0.567	0.523	0.561
Video	MAE	1.422	-	1.441	1.532	1.419
	Corr	0.138	-	0.205	0.119	0.156
	F1	0.535	-	0.552	0.534	0.513
	Acc	0.512	-	0.542	0.421	0.512
Full	MAE	1.191	-	1.093	1.055	1.011
	Corr	0.463	-	0.612	0.607	0.663
	F1	0.726	-	0.735	0.757	0.764
	Acc	0.696	-	0.736	0.763	0.765
Text	AUROC	0.641	0.527	0.67	0.654	0.67
	Accuracy	0.585	0.522	0.607	0.515	0.63
Image	AUROC	0.639	0.541	0.64	0.653	0.665
	Accuracy	0.595	0.517	0.591	0.515	0.622
Full	AUROC	0.649	0.683	0.633	0.666	0.675
	Accuracy	0.585	0.634	0.596	0.536	0.634

Table 8. Supervised learning results on CMU-MOSI and Hateful Memes datasets.

25% percentile within the batch) effectively filters out some false positives identified by the pseudo labels.

D.9 Clusterability Analysis

Complementing the Alignment Uniformity analysis presented in the main manuscript, we provide a comparison of the clusterability characteristic of the learned representations in Figure 10. This experiment is conducted with Robult, compared to GMC and ActionMAE on Hateful Memes dataset. This qualitative analysis demonstrates that Robult’s representations are better clustered even in scenarios where different modalities are missing. In contrast, the other methods do not exhibit this level of clustering effectiveness.

Modality	Metrics	MOSI Dataset			MOSEI Dataset		
		DiCMoR	SEM	Robult	DiCMoR	SEM	Robult
Text	MAE	1.444	1.632	1.397	0.819	0.894	0.784
	Corr	0.085	0.095	0.144	0.276	0.252	0.459
	F1	0.536	0.506	0.578	0.572	0.599	0.739
	Acc	0.511	0.53	0.569	0.657	0.63	0.732
Audio	MAE	1.504	1.739	1.415	0.829	1.037	0.825
	Corr	0.006	0.043	0.085	0.201	0.138	0.221
	F1	0.484	0.441	0.539	0.537	0.545	0.679
	Acc	0.49	0.463	0.535	0.642	0.536	0.65
Vision	MAE	1.454	1.87	1.425	0.83	0.992	0.826
	Corr	0.019	0.017	0.086	0.163	0.133	0.201
	F1	0.526	0.449	0.593	0.552	0.523	0.647
	Acc	0.524	0.475	0.522	0.635	0.524	0.632
Text + Audio	MAE	1.481	1.728	1.401	0.828	0.919	0.762
	Corr	0.013	0.125	0.141	0.173	0.248	0.439
	F1	0.494	0.443	0.563	0.563	0.611	0.733
	Acc	0.495	0.465	0.546	0.639	0.602	0.717
Text + Vision	MAE	1.473	1.758	1.415	0.832	0.92	0.788
	Corr	0.012	0.077	0.146	0.158	0.25	0.399
	F1	0.514	0.452	0.58	0.579	0.612	0.718
	Acc	0.514	0.476	0.534	0.638	0.592	0.704
Audio + Vision	MAE	1.478	1.794	1.426	0.836	0.923	0.812
	Corr	0.023	0.035	0.091	0.138	0.143	0.244
	F1	0.484	0.429	0.581	0.581	0.535	0.663
	Acc	0.491	0.451	0.527	0.626	0.552	0.64
Full	MAE	1.468	1.797	1.392	0.839	0.902	0.779
	Corr	0.035	0.041	0.247	0.149	0.249	0.504
	F1	0.488	0.432	0.657	0.587	0.625	0.744
	Acc	0.495	0.453	0.63	0.614	0.667	0.741

Table 9. Additional comparison on CMU-MOSI and CMU-MOSEI Datasets.

D.10 Transferability Analysis

With this experiment, we investigate the tranferability characteristic of Robult, as well as existing state-of-the-art frameworks and baselines.

Experiment settings. Inspired by common pre-training procedures, where a model is initially trained on a large dataset for a source task and then fine-tuned for a target task, we designed an experiment to evaluate the zero-shot performance of all models on CMU-MOSI after being trained with the CMU-MOSEI dataset. This setting aligns with common practices, as CMU-MOSEI is larger in scale, covering a wider range of sentiment levels and emotions compared to CMU-MOSI [Zadeh *et al.*, 2018b]. To conduct the experiment, we first pretrain all methods with CMU-MOSEI using 5% labeled data, simultaneously evaluating them in a semi-supervised scenario. Since zero-shot evaluation requires no fine-tuning stage, all model architectures must remain intact after pretraining. However, there are discrepancies in the dimensions of the input data between CMU-MOSI and CMU-MOSEI. Specifically, the audio input of CMU-MOSI has a latent dimension of 5, while that of CMU-MOSEI is 74. Additionally, the video input of CMU-MOSI and CMU-MOSEI is 20 and 35, respectively. To address this issue, we generate a compact version of CMU-MOSEI by employing T-SNE on the original data, aligning the dimensions with those in the CMU-MOSI datasets. After pretraining, we directly run evaluations on the normal test set of CMU-MOSI, given different

Modality	Metrics	Framework		
		Robult	Robult w/o weighting scheme	Robult w/o unique branches
<i>MOSI Dataset:</i>				
Text	MAE	1.397	1.418	1.514
	Corr	0.144	0.125	0.131
	F1	0.578	0.551	0.53
	Acc	0.569	0.548	0.443
Audio	MAE	1.415	1.479	1.576
	Corr	0.085	-0.042	-0.096
	F1	0.539	0.514	0.513
	Acc	0.535	0.456	0.514
Vision	MAE	1.425	1.434	1.509
	Corr	0.086	0.087	0.034
	F1	0.593	0.593	0.526
	Acc	0.522	0.422	0.528
Full	MAE	1.392	1.388	1.487
	Corr	0.247	0.192	0.207
	F1	0.657	0.566	0.567
	Acc	0.63	0.569	0.496
<i>Hateful Memes:</i>				
Text	AUROC	0.623	0.556	0.586
	Accuracy	0.59	0.541	0.577
Image	AUROC	0.596	0.597	0.547
	Accuracy	0.562	0.511	0.533
Full	AUROC	0.632	0.571	0.601
	Accuracy	0.595	0.345	0.544

Table 10. Additional Ablation Study with Robult on two datasets CMU-MOSI and Hateful Memes.

Modality	Metrics	MOSI Dataset		MOSEI Dataset	
		Robult - L1	Robult - L2	Robult	Robult - L1
Text	MAE	1.486	1.456	1.397	0.793
	Corr	0.1	0.184	0.144	0.421
	F1	0.571	0.573	0.578	0.733
	Acc	0.545	0.576	0.569	0.729
Audio	MAE	1.475	1.51	1.415	0.825
	Corr	0.049	0.083	0.085	0.199
	F1	0.544	0.477	0.539	0.674
	Acc	0.52	0.478	0.535	0.635
Vision	MAE	1.475	1.478	1.425	0.917
	Corr	0.028	0.045	0.086	0.133
	F1	0.593	0.582	0.593	0.567
	Acc	0.492	0.522	0.522	0.572
Text + Audio	MAE	1.477	1.395	1.401	0.764
	Corr	0.089	0.166	0.141	0.439
	F1	0.57	0.558	0.563	0.74
	Acc	0.531	0.561	0.546	0.722
Text + Vision	MAE	1.47	1.389	1.415	0.781
	Corr	0.128	0.236	0.146	0.391
	F1	0.59	0.553	0.58	0.705
	Acc	0.52	0.556	0.534	0.7
Audio + Vision	MAE	1.465	1.475	1.426	0.834
	Corr	0.04	0.054	0.091	0.187
	F1	0.577	0.533	0.581	0.635
	Acc	0.472	0.491	0.527	0.622
Full	MAE	1.403	1.366	1.392	0.812
	Corr	0.223	0.235	0.247	0.45
	F1	0.554	0.585	0.657	0.703
	Acc	0.547	0.583	0.63	0.708

Table 11. Additional ablation study on CMU-MOSI and CMU-MOSEI Datasets.

combinations of input modalities to evaluate modalities missing performance.

Results. Table 13 provides a summary of the results from this

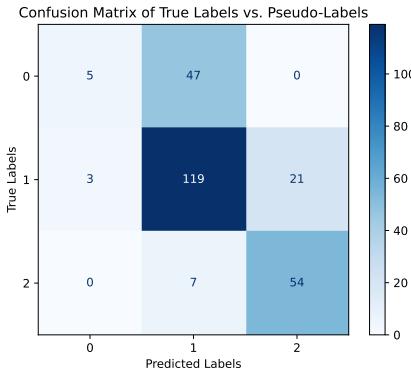
Modality	Mutual Information with fused representation	
	Robult	Robult w/o Soft P-U Loss
Text	0.309	0.054
Audio	0.285	0.077
Vision	0.274	0.083
Fused	2.037	1.707

Table 12. Mutual Information between fused and unimodal representations on the CMU-MOSI dataset.

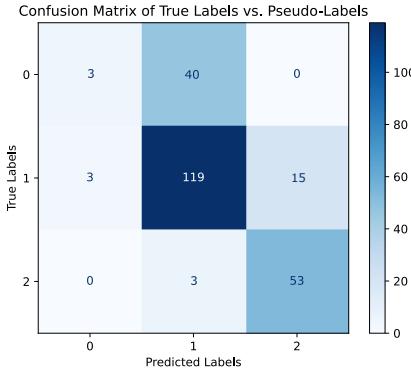
Metrics	MOSI Dataset			
	Unimodal	GMC	ActionMAE	Robult
<i>Text Modality:</i>				
MAE		1.448	1.454	1.456
Corr		0.132	0.119	0.106
F1		0.559	0.551	0.624
Acc		0.527	0.48	0.481
<i>Audio Modality:</i>				
MAE		1.514	1.456	1.517
Corr		-0.196	0.085	0.084
F1		0.5	0.592	0.57
Acc		0.512	0.429	0.423
<i>Vision Modality:</i>				
MAE		1.473	1.86	1.547
Corr		0.076	-0.087	0.017
F1		0.71	0.593	0.561
Acc		0.546	0.422	0.432
<i>Text+Audio Modalities:</i>				
MAE		1.405	1.441	1.378
Corr		0.013	0.164	0.101
F1		0.504	0.569	0.591
Acc		0.456	0.467	0.49
<i>Text+Vision Modalities:</i>				
MAE		1.45	1.629	1.403
Corr		0.119	0.106	0.137
F1		0.621	0.593	0.606
Acc		0.536	0.422	0.568
<i>Audio+Vision Modalities:</i>				
MAE		1.434	1.57	1.528
Corr		-0.181	-0.005	0.08
F1		0.496	0.533	0.551
Acc		0.459	0.422	0.441
<i>Full Modalities:</i>				
MAE		1.456	1.684	1.472
Corr		0.088	0.072	0.066
F1		0.549	0.573	0.55
Acc		0.551	0.425	0.551

Table 13. Transferability result on CMU-MOSI dataset.

experiment. Generally, all methods experience a reduction in performance in certain cases when transferred to a different dataset. However, among all approaches, Robult consistently achieves the best performance, as indicated by the recorded metrics. In addition, it is noteworthy that Robult is the only approach capable of producing meaningful results with input



(a) Without weighting scheme



(b) With weighting scheme: Weight threshold of 25% percentile.

Figure 9. Confusion matrix of Pseudo Labels versus groundtruth label at epoch 20 on CMU-MOSI dataset.

from full modalities in this zero-shot transfer setting.

D.11 Incorporation with existing approaches

This analysis investigates the ability of Robult in incorporating with other approaches to enhance their desired characteristics in learned representations.

Experiment settings. We select GMC as a baseline approach for conducting this experiment. GMC aims to preserve the geometrical alignment of representations from different modalities through a geometrical contrastive loss [Poklukar *et al.*, 2022b]. To observe the impact of incorporating Robult with GMC to preserve this characteristic, we simply adopt their geometrical contrastive loss with our existing

$$\mathcal{L}_{(u)lb}:$$

$$\begin{aligned} \mathcal{L}_{lb}^i &= -\frac{1}{\|B_{1,1}\|} \sum_{\substack{(j,k) \sim \\ p(F=1, L=1)}} \log v(s_j, z_k^i) + \log v(s_j, s_k) \\ &\quad + \log v(z_j^i, z_k^i); \\ \mathcal{L}_{lb} &= -\frac{1}{M} \sum_{i=1}^M \mathcal{L}_{lb}^i. \end{aligned} \quad (18)$$

and:

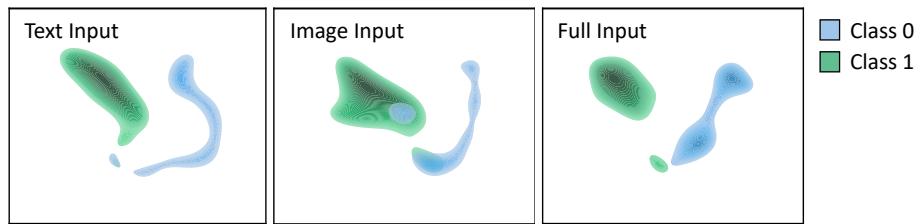
$$\begin{aligned} \mathcal{L}_{ulb}^i &= -\frac{1}{\|B_{1,0}\|} \sum_{\substack{(j,k) \sim \\ p(F=1, L=0)}} w_{jk}^i (\log v(s_j, z_k^i) + \log v(s_j, s_k) \\ &\quad + \log v(z_j^i, z_k^i)); \\ \mathcal{L}_{ulb} &= -\frac{1}{M} \sum_{i=1}^M \mathcal{L}_{ulb}^i. \end{aligned} \quad (19)$$

To evaluate the geometrical alignment of the learned representations, we employ Delaunay Component Analysis (DCA) [Poklukar *et al.*, 2022a], a technique similar to that used in GMC. DCA involves comparing geometric and topological properties of an evaluation set of representations (E) with a reference set (R), which acts as an approximation of the true underlying manifold. Following the evaluation strategy outlined in [Poklukar *et al.*, 2022b], we consider three metrics provided by DCA that reflect the geometric alignment between R (representations of full modalities input) and E (representations of single modality inputs): network quality $q \in [0, 1]$, precision \mathcal{P} , and recall \mathcal{R} . We report the harmonic mean defined as $3/(1/\mathcal{P} + 1/\mathcal{R} + 1/q)$ when all $\mathcal{P}, \mathcal{R}, q > 0$ and 0 otherwise. For a detailed description of DCA and its settings, please refer to the original work [Poklukar *et al.*, 2022a; Poklukar *et al.*, 2022b].

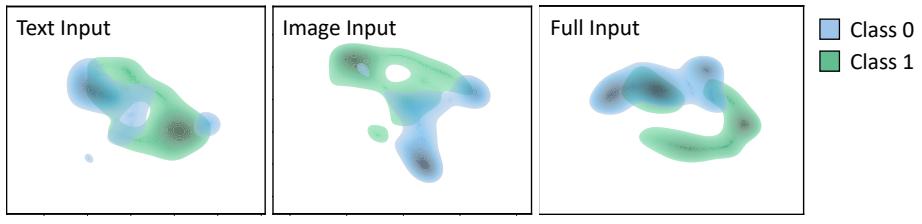
Results. We provide the alignment metrics for the representations generated with CMU-MOSI and Hateful Memes datasets, considering only 50% labeled data in their respective training sets (Table 14). The statistics indicate that Robult effectively enhances the performance of GMC in its effort to preserve geometrical alignment under the constraint of limited label information. We anticipate that this behavior can potentially be extended to other methods under limited available label information, although additional investigations are needed to verify this.

Dataset	R	E	Metrics	Unimodal	GMC	Robult + GMC
MOSI Dataset	Full	Text	MAE	0.473	0.529	0.535
	Full	Audio	Corr	0	0.375	0.393
	Full	Vision	F1	0	0.478	0.335
Hateful Memes	Full	Text	AUROC	0.349	0.489	0.518
	Full	Image	AUROC	0	0.456	0.509

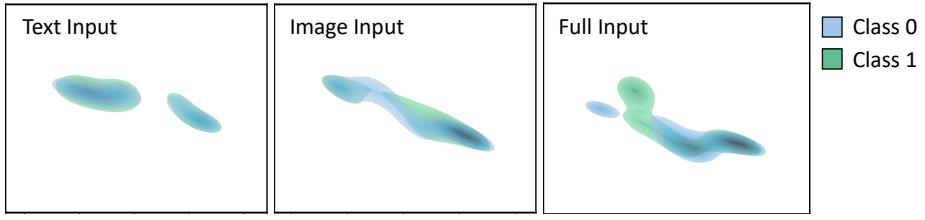
Table 14. DCA Scores of models, evaluating geometrical alignment of full-modalities representations with unimodal representations.



(a) Robust



(b) GMC



(c) ActionMAE

Figure 10. Representation clusters generated by different methods on Hateful Memes dataset.