M3CoL: HARNESSING SHARED RELATIONS VIA MULTIMODAL MIXUP CONTRASTIVE LEARNING FOR MULTIMODAL CLASSIFICATION

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028 029

031

Paper under double-blind review

ABSTRACT

Deep multimodal learning has shown remarkable success by leveraging contrastive learning to capture explicit one-to-one relations across modalities. However, realworld data often exhibits shared relations beyond simple pairwise associations. We propose M3CoL, a Multimodal Mixup Contrastive Learning approach to capture nuanced shared relations inherent in multimodal data. Our key contribution is a Mixup-based contrastive loss that learns robust representations by aligning mixed samples from one modality with their corresponding samples from other modalities thereby capturing shared relations between them. For multimodal classification tasks, we introduce a framework that integrates a fusion module with unimodal prediction modules for auxiliary supervision during training, complemented by our proposed Mixup-based contrastive loss. Through extensive experiments on diverse datasets (N24News, ROSMAP, BRCA, and Food-101), we demonstrate that M3CoL effectively captures shared multimodal relations and generalizes across domains. It outperforms state-of-the-art methods on N24News, ROSMAP, and BRCA, while achieving comparable performance on Food-101. Our work highlights the significance of learning shared relations for robust multimodal learning, opening up promising avenues for future research.

1 INTRODUCTION

The way we perceive the world is shaped by various modalities, such as language, vision, audio, 033 and more. In the era of abundant and accessible multimodal data, it is increasingly crucial to 034 equip artificial intelligence with multimodal capabilities (Baltrušaitis et al., 2018). At the heart of advancements in multimodal learning is contrastive learning, which maximizes similarity for positive pairs and minimizes it for negative pairs, making it practical for multimodal representation learning. CLIP (Radford et al., 2021) is a prominent example that employs contrastive learning to understand 037 the direct link between paired modalities and seamlessly maps images and text into a shared space for cross-modal understanding, which can be later utilized for tasks such as retrieval and classification. However, traditional contrastive learning methods often overlook shared relationships between 040 samples across different modalities, which can result in the learning of representations that are not 041 fully optimized for capturing the underlying connections between diverse data modalities. These 042 methods primarily focus on distinguishing between positive and negative pairs of samples, typically 043 treating each instance as an independent entity. They tend to disregard the rich, shared relational 044 information that could exist between samples within and across modalities. This limited focus can prevent the model from leveraging valuable contextual information, such as semantic similarities or complementary patterns, which can enhance robust representation learning. Consequently, this can 046 lead to suboptimal performance in downstream tasks that require optimized shared representations, 047 such as image-text alignment, cross-modal retrieval, or multimodal fusion tasks. 048

As shown in the left panel of Figure 1, classical contrastive learning approach assumes perfect
 one-to-one relations between modalities, which is rare in real-world data. For example, shared
 elements in images or text can relate even across separate samples, as illustrated by the elements
 like "tomato sauce" and "basil" in Figure 1. Our approach, illustrated in the right panel of Figure 1,
 goes beyond simple pairwise alignment by capturing shared relationships across mixed samples. By
 creating newer data points through convex combinations of data points our method more effectively

056

058

060

061

062

063

064

065

066 067

068

069

071

096

098

099

102

103



Figure 1: Comparison of traditional contrastive and our proposed M3Co loss. $\mathbf{M}_{i}^{(1)}$ and $\mathbf{M}_{i}^{(2)}$ denote representations of the *i*-th sample from modalities 1 and 2, respectively. Traditional contrastive loss (left panel) aligns corresponding sample representations across modalities. M3Co (right panel) mixes the *i*-th and *j*-th samples from modality 1 and enforces the representations of this mixture to align with the representations of the corresponding *i*-th and *j*-th samples from modality 2, and vice versa. For the text modality, we mix the text embeddings, while we mix the raw inputs for other modalities. 073 Similarity (Sim) represents type of alignment enforced between the embeddings for all modalities. 074

075 models complex shared relationships, such as imperfect bijections (Liang et al., 2022a), enhancing 076 multimodal classification performance.

077 Our approach builds upon the success of data augmentation techniques such as Mixup (Zhang et al., 078 2017) and their variants (Yun et al., 2019; Cubuk et al., 2019; Hendrycks et al., 2019), which have 079 proven beneficial for enhancing learned feature spaces, improving both robustness and performance. Mixup trains models on synthetic data created through convex combinations of two datapoint-label 081 pairs (Chapelle et al., 2000). These techniques are particularly valuable in low sample settings, as 082 they help prevent overfitting and the learning of ineffective shortcuts (Chen et al., 2020; Robinson et al., 2021), common in contrastive learning. Building on the success of recent Mixup strategies (Shen et al., 2022; Thulasidasan et al., 2019; Verma et al., 2019) and MixCo (Kim et al., 2020), 084 we introduce M3Co, a novel approach that significantly adapts and enhances contrastive learning 085 principles to complex multimodal settings. M3Co modifies the CLIP loss to effectively handle multimodal scenarios, addressing the problem of instance discrimination, where models overly 087 focus on distinguishing individual instances instead of capturing relationships between modalities. By leveraging convex combinations of data for contrastive learning, M3Co eliminates instance discrimination and enhances robust representation learning by capturing shared relations. These 090 combinations serve as structured noise and treated as positive pairs with their corresponding samples 091 from other modalities. Our experimental results demonstrate enhanced ability to capture shared 092 relations enabling improvements in performance and generalization across a range of multimodal classification tasks.

- 094 Our key contributions are summarized as follows: 095
 - We propose M3Co, a multimodal contrastive loss (Eq. 8) that utilizes mixed samples to effectively capture shared relationships across different modalities. By going beyond traditional pairwise alignment methods, M3Co makes representations more consistent with the complex, intertwined relationships usually observed in real-world data.
 - · We introduce a multimodal learning framework (Figure 2) consisting of unimodal prediction modules, a fusion module, and a novel Mixup-based contrastive loss. Our proposed method is modality-agnostic, allowing for flexible application across various types of data, and continuously updates the representations necessary for accurate and consistent predictions.
- We demonstrate the effectiveness of our methodology by evaluating it on four public 105 multimodal classification benchmark datasets from different domains: two image-text 106 datasets, N24News and Food-101, and two medical datasets, ROSMAP and BRCA (Table 1, 107 2, 3). Our approach outperforms baseline models, especially on smaller datasets.

¹⁰⁸ 2 METHODOLOGY

109 110 111

112

113

114

115

116

117

118

119

120

121 122 123

124

125

127 128 129

130 131 132

133 134

135 136

137 138

148

149

Pipeline Overview: Figure 2 depicts our framework, which comprises of three components: unimodal prediction modules, a fusion module, and a Mixup-based contrastive loss. We obtain latent representations (using learnable modality specific encoders $f^{(1)}$ and $f^{(2)}$) of individual modalities and fuse them (denoted by concatenation symbol '+') to generate a joint multimodal representation, which is optimized using a supervised objective (through classifier 3). The unimodal prediction modules provide additional supervision during training (via classifier 1 and 2). These strategies enable deeper integration of modalities and allow the models to compensate for the weaknesses of one modality with the strengths of another. The Mixup-based contrastive loss (denoted by \mathcal{L}_{M3Co}) continuously updates the representations by capturing shared relations inherent in the multimodal data. This comprehensive approach enhances the understanding of multimodal data, improving accuracy and model robustness.



Figure 2: Architecture of our proposed M3CoL model. Samples from modality $1 (\mathbf{x}_i^{(1)}, \mathbf{x}_j^{(1)})$ and modality $2 (\mathbf{x}_i^{(2)}, \mathbf{x}_k^{(2)})$, along with their respective mixed data $\tilde{\mathbf{x}}_{i,j}^{(1)}$ and $\tilde{\mathbf{x}}_{i,k}^{(2)}$, are fed into encoders $f^{(1)}$ and $f^{(2)}$ to generate embeddings. Unimodal embeddings $\mathbf{p}_i^{(1)}$ and $\mathbf{p}_i^{(2)}$ are processed through classifier 1 and 2 to produce predictions $\hat{\mathbf{y}}_i^{(1)}$ and $\hat{\mathbf{y}}_i^{(2)}$ for training supervision only. The unimodal embeddings $\mathbf{p}_i^{(1)}$ and $\mathbf{p}_i^{(2)}$ are concatenated and processed through classifier 3 to yield $\hat{\mathbf{y}}_{\text{final}}$, utilized during training and inference. Additionally, unimodal embeddings $\mathbf{p}_i^{(1)}, \mathbf{p}_j^{(1)}, \mathbf{p}_i^{(2)}, \mathbf{p}_k^{(2)}$, and mixed embeddings $\tilde{\mathbf{p}}_{i,j}^{(1)}$ and $\tilde{\mathbf{p}}_{i,k}^{(2)}$ are utilized by our contrastive loss $\mathcal{L}_{\text{M3Co}}$ for shared alignment.

Multimodal Mixup Contrastive Learning: Given a batch of N multimodal samples, let $\mathbf{x}_i^{(1)}$ and $\mathbf{x}_i^{(2)}$ denote the *i*-th samples for the first and second modalities, respectively. The modality encoders, $f^{(1)}$ and $f^{(2)}$, generate the corresponding embeddings $\mathbf{p}_i^{(1)}$ and $\mathbf{p}_i^{(2)}$:

$$\mathbf{p}_{i}^{(1)} = f^{(1)}(\mathbf{x}_{i}^{(1)}), \quad \mathbf{p}_{i}^{(2)} = f^{(2)}(\mathbf{x}_{i}^{(2)})$$
(1)

155 156 We generate a mixture, $\tilde{\mathbf{x}}_{i,j}^{(1)}$, of the samples $\mathbf{x}_i^{(1)}$ and $\mathbf{x}_j^{(1)}$ by taking their convex combination. 157 Similarly, we generate a mixture, $\tilde{\mathbf{x}}_{i,k}^{(2)}$, using the convex combination of the samples $\mathbf{x}_i^{(2)}$ and $\mathbf{x}_k^{(2)}$ 158 (Eq. 2). In the case of text modality, instead of directly mixing the raw inputs, we mix the text 159 embeddings (Guo et al., 2019). The mixing indices j, k are drawn arbitrarily, without replacement, 160 from [1, N], for both the modalities. We mix both the modalities using a factor $\lambda \sim \text{Beta}(\alpha, \alpha)$. 161 Based on the findings of (Zhang et al., 2017), which demonstrated enhanced performance for α values 162 between 0.1 and 0.4, we chose $\alpha = 0.15$ after experimenting with several values in this range. The

mixtures are fed through the respective encoders to obtain the embeddings: $\tilde{\mathbf{p}}_{i,j}^{(1)}$, and $\tilde{\mathbf{p}}_{i,k}^{(2)}$ (Eq. 3).

$$\tilde{\mathbf{x}}_{i,j}^{(1)} = \lambda_i \cdot \mathbf{x}_i^{(1)} + (1 - \lambda_i) \cdot \mathbf{x}_j^{(1)}, \quad \tilde{\mathbf{x}}_{i,k}^{(2)} = \lambda_i \cdot \mathbf{x}_i^{(2)} + (1 - \lambda_i) \cdot \mathbf{x}_k^{(2)}$$
(2)

$$\tilde{\mathbf{p}}_{i}^{(1)} = \tilde{\mathbf{p}}_{i,j}^{(1)} = f^{(1)}(\tilde{\mathbf{x}}_{i,j}^{(1)}), \quad \tilde{\mathbf{p}}_{i}^{(2)} = \tilde{\mathbf{p}}_{i,k}^{(2)} = f^{(2)}(\tilde{\mathbf{x}}_{i,k}^{(2)})$$
(3)

The unidirectional contrastive loss (Sohn, 2016; Chen et al., 2020; Oord et al., 2018; Wu et al., 2018; Zhang et al., 2022) over $\mathbf{p}^{(2)}$ is conventionally defined as:

$$\mathcal{L}_{\text{sim-conv}}(\mathbf{p}^{(1)}, \mathbf{p}^{(2)}) = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp\left(\mathbf{p}_{i}^{(1)} \cdot \mathbf{p}_{i}^{(2)} / \tau\right)}{\sum_{j=1}^{N} \exp\left(\mathbf{p}_{i}^{(1)} \cdot \mathbf{p}_{j}^{(2)} / \tau\right)}$$
(4)

where \cdot indicates dot product and τ is a temperature hyperparameter. While this formulation is needed for computing similarity among aligned samples from different modalities, our loss handles both aligned and non-aligned samples, as this enables to learn a better representation space. To achieve this, we define the unidirectional multimodal contrastive loss between $\mathbf{p}_i^{(1)}$ and $\mathbf{p}_m^{(2)}$ over $\mathbf{p}^{(2)}$ as:

$$\mathcal{L}_{\rm sim}(\mathbf{p}_i^{(1)}, \mathbf{p}^{(2)}; m) = -\log \frac{\exp\left(\mathbf{p}_i^{(1)} \cdot \mathbf{p}_m^{(2)} / \tau\right)}{\sum\limits_{j=1}^N \exp\left(\mathbf{p}_i^{(1)} \cdot \mathbf{p}_j^{(2)} / \tau\right)}$$
(5)

where $\mathbf{p}^{(1)}$ and $\mathbf{p}^{(2)}$ are \mathcal{L}^2 normalized, τ is a temperature hyperparameter, and *m* is a sample index in [1, *N*]. Although the unidirectional multimodal contrastive loss (Eq. 5) can learn indirect relations, it is insufficient for learning shared semi-positive relations between modalities. Therefore, we introduce a Mixup-based contrastive loss to capture these relations that promotes generalized learning, as this process is more nuanced than simply discriminating positives from negatives. Now, we make our loss bidirectional to encourage improved alignment in the shared representation space and efficient use of training data (Radford et al., 2021; Oord et al., 2018; Sohn, 2016). We define this bidirectional Mixup contrastive loss M3Co for each modality (Eq. 6, 7) and the total M3Co loss (Eq. 8) as:

 $\mathcal{L}_{\text{M3Co}}^{(1)} = \frac{1}{N} \sum_{i=1}^{N} \left[\lambda_i \cdot \mathcal{L}_{\text{sim}}(\tilde{\mathbf{p}}_{i,j}^{(1)}, \mathbf{p}^{(2)}; i) + (1 - \lambda_i) \cdot \mathcal{L}_{\text{sim}}(\tilde{\mathbf{p}}_{i,j}^{(1)}, \mathbf{p}^{(2)}; j) \right]$

 $+\frac{1}{N}\sum_{i=1}^{N}\left\{\lambda_{i}\cdot\mathcal{L}_{\text{sim}}(\mathbf{p}_{i}^{(2)},\tilde{\mathbf{p}}^{(1)};i)+(1-\lambda_{i})\cdot\mathcal{L}_{\text{sim}}(\mathbf{p}_{j}^{(2)},\tilde{\mathbf{p}}^{(1)};i)\right\}$

$$\mathcal{L}_{\text{M3Co}}^{(2)} = \frac{1}{N} \sum_{i=1}^{N} \left[\lambda_i \cdot \mathcal{L}_{\text{sim}}(\tilde{\mathbf{p}}_{i,k}^{(2)}, \mathbf{p}^{(1)}; i) + (1 - \lambda_i) \cdot \mathcal{L}_{\text{sim}}(\tilde{\mathbf{p}}_{i,k}^{(2)}, \mathbf{p}^{(1)}; k) \right] \\ + \frac{1}{N} \sum_{i=1}^{N} \left\{ \lambda_i \cdot \mathcal{L}_{\text{sim}}(\mathbf{p}_i^{(1)}, \tilde{\mathbf{p}}^{(2)}; i) + (1 - \lambda_i) \cdot \mathcal{L}_{\text{sim}}(\mathbf{p}_k^{(1)}, \tilde{\mathbf{p}}^{(2)}; i) \right\}$$

$$\mathcal{L}_{\rm M3Co}^{(1,2)} = \frac{1}{2} \left(\mathcal{L}_{\rm M3Co}^{(1)} + \mathcal{L}_{\rm M3Co}^{(2)} \right) \tag{8}$$

(6)

(7)

where $\mathbf{p}^{(1)}$, $\tilde{\mathbf{p}}^{(1)}$, $\mathbf{p}^{(2)}$, and $\tilde{\mathbf{p}}^{(2)}$ are \mathcal{L}^2 normalized. Note that the parts of the loss functions in Eq. (6, 7) inside curly parantheses make them bidirectional. Mixup-based methods enhance generalization by capturing clean patterns in the early training stages but can eventually overfit to noise if continued for larger number of epochs (Liu et al., 2023; Yu et al., 2021; Golatkar et al., 2019). To address this, we implement a schedule that transitions from the Mixup-based M3Co loss to a non-Mixup multimodal contrastive loss. We design this transition so that the non-Mixup loss retains the ability to learn shared or indirect relationships between modalities. By using a bidirectional SoftClip-based loss (Gao et al., 2024; Sohn, 2016; Chen et al., 2020), we relax the rigid one-to-one correspondence, allowing the model to capture many-to-many relations (Gao et al., 2024; 2022). The bidirectional

217218219220221

222

224 225 226

227 228 229

243

244 245

250

256 257 258

259

MultiSoftClip loss for each modality (Eq. 9, 10) and its combination (Eq. 11) is:

$$\mathcal{L}_{\text{MultiSClip}}^{(1)} = \frac{1}{N} \sum_{i=1}^{N} \sum_{l=1}^{N} \left[\frac{\exp\left(\mathbf{p}_{i}^{(1)} \cdot \mathbf{p}_{l}^{(1)} / \tau\right)}{\sum_{t=1}^{N} \exp\left(\mathbf{p}_{i}^{(1)} \cdot \mathbf{p}_{t}^{(1)} / \tau\right)} \cdot \left(\mathcal{L}_{\text{sim}}(\mathbf{p}_{i}^{(2)}, \mathbf{p}^{(1)}; l) + \mathcal{L}_{\text{sim}}(\mathbf{p}_{l}^{(1)}, \mathbf{p}^{(2)}; i) \right) \right]$$
(9)

$$\mathcal{L}_{\text{MultiSClip}}^{(2)} = \frac{1}{N} \sum_{i=1}^{N} \sum_{l=1}^{N} \left[\frac{\exp\left(\mathbf{p}_{i}^{(2)} \cdot \mathbf{p}_{l}^{(2)} / \tau\right)}{\sum_{t=1}^{N} \exp\left(\mathbf{p}_{i}^{(2)} \cdot \mathbf{p}_{t}^{(2)} / \tau\right)} \cdot \left(\mathcal{L}_{\text{sim}}(\mathbf{p}_{i}^{(1)}, \mathbf{p}^{(2)}; l) + \mathcal{L}_{\text{sim}}(\mathbf{p}_{l}^{(2)}, \mathbf{p}^{(1)}; i) \right) \right]$$
(10)

$$\mathcal{L}_{\text{MultiSClip}}^{(1,2)} = \frac{1}{2} \left(\mathcal{L}_{\text{MultiSClip}}^{(1)} + \mathcal{L}_{\text{MultiSClip}}^{(2)} \right)$$
(11)

where $\mathbf{p}^{(1)}$ and $\mathbf{p}^{(2)}$ are \mathcal{L}^2 normalized. The M3Co and MultiSClip losses for M modalities is:

$$\mathcal{L}_{\text{M3Co}} = \sum_{i=1}^{M} \sum_{j>i}^{M} \mathcal{L}_{\text{M3Co}}^{(i,j)}$$
(12)

$$\mathcal{L}_{\text{MultiSClip}} = \sum_{i=1}^{M} \sum_{j>i}^{M} \mathcal{L}_{\text{MultiSClip}}^{(i,j)}$$
(13)

Unimodal Predictions and Fusion: The encoders produce latent representations for each of the Mmodalities, serving as inputs to individual classifiers that generate modality-specific predictions $\hat{\mathbf{y}}^{(m)}$. These representations are used for modality-specific supervision only during training. The unimodal prediction task involves minimizing the cross-entropy loss \mathcal{L}_{CE} between these predictions and the corresponding ground truth labels (\mathbf{y}), for each modality. The unimodal cross-entropy loss is:

$$\mathcal{L}_{\text{CE-Uni}} = \sum_{m=1}^{M} \mathcal{L}_{\text{CE}}(\mathbf{y}, \hat{\mathbf{y}}^{(m)})$$
(14)

We merge the unimodal latent representations by concatenating them and pass the combined representation to the output classifier. These predictions serve as the final outputs $\hat{\mathbf{y}}_f$ used during inference. The multimodal prediction process aims to minimize the cross-entropy loss between $\hat{\mathbf{y}}_f$ and the corresponding labels. The multimodal cross-entropy loss is:

$$\mathcal{L}_{\text{CE-Multi}} = \mathcal{L}_{\text{CE}}(\mathbf{y}, \hat{\mathbf{y}}_f) \tag{15}$$

Combined Learning Objective: Our overall loss objective utilizes a schedule to combine our M3Co and MultiSClip loss functions weighted by a hyperparamater β , along with the unimodal and multimodal cross-entropy losses. We use M3Co for the first one-third (Liu et al., 2023) part of training, and then transition to MultiSClip as over-training with a Mixup-based loss can potentially harm generalization. The end-to-end loss is defined as:

$$\mathcal{L}_{\text{Total}} = \beta \cdot \mathcal{L}_{\text{M3Co} \mid \text{MultiSClip}} + \mathcal{L}_{\text{CE-Uni}} + \mathcal{L}_{\text{CE-Multi}}$$
(16)

3 EXPERIMENTS

260 **Datasets.** We evaluate our approach on four diverse publicly available multimodal classification 261 datasets: N24News (Wang et al., 2022), Food-101 (Wang et al., 2015), ROSMAP (Wang et al., 262 2021a), and BRCA (Wang et al., 2021a). N24News and Food-101 are both bimodal image-text classification datasets. Food-101 is a food classification dataset, where each sample is linked 264 with a recipe description gathered from web pages and an associated image. N24News is a news 265 classification dataset consisting of four text types (Abstract, Caption, Heading, and Body) along 266 with the corresponding images. Following other works (Zou et al., 2023), we use the first three text types for our experiments. ROSMAP and BRCA are publicly available multimodal medical datasets, 267 each containing three modalities: DNA methylation, miRNA expression, and mRNA expression. 268 **ROSMAP** is an Alzeihmer's diagnosis dataset, while **BRCA** is used for breast invasive carcinoma 269 PAM50 subtype classification. Appendix A.2 provides information about the train-val-test splits.

Evaluation Metrics. The evaluation metric used for N24News and Food-101 is classification
accuracy (ACC). For BRCA, we report accuracy (ACC), macro-averaged F1 score (MF1), and
weighted F1 score (WF1). For ROSMAP, we use accuracy (ACC), area under the ROC curve (AUC),
and F1 score (F1) as the evaluation metrics.

Implementation Details. We use a ViT (pre-trained on the ImageNet-21k dataset) (Dosovitskiy et al., 2020) as the image encoder for N24News and Food-101. For N24News, the text encoder is a pretrained BERT/RoBERTa (Devlin et al., 2018; Zhuang et al., 2021), while we use a pretrained BERT as the text encoder for Food-101. The classifiers for the above two datasets are three layer MLPs with ReLU activations. For ROSMAP and BRCA, which are small datasets, we use two layer MLPs as feature encoders for each modality, and two layer MLPs with ReLU activations as classifiers. The hyperparameter settings and all other details are given in Appendix A.1.

281 **Baselines.** We compare our method with various multimodal classification approaches (Van De Wiel 282 et al., 2016; Wang et al., 2021a; Han et al., 2020; Abavisani et al., 2020; Han et al., 2022; Zou et al., 283 2023; Liang et al., 2022b; Kiela et al., 2019; 2018; Wang et al., 2022; Vielzeuf et al., 2018; Arevalo 284 et al., 2017; Li et al., 2019; Huang et al., 2020; Kim et al., 2021; Narayana et al., 2019; Liu et al., 285 2021; Hong et al., 2020; Huang et al., 2021; Singh et al., 2019; Wang et al., 2024). Some methods (Kiela et al., 2019; Vielzeuf et al., 2018; Arevalo et al., 2017) focus on integrating global features 286 287 from individual modality-specific backbones to enhance classification. Others (Li et al., 2019; Kim et al., 2021; Huang et al., 2020; Narayana et al., 2019) use sophisticated pre-trained architectures 288 fine-tuned for specific tasks. UniS-MMC (Zou et al., 2023), the previous state-of-the-art on Food-289 101 and N24News, uses contrastive learning to align features across modalities with supervision 290 from unimodal predictions. Similarly, Dynamics (Han et al., 2022), the previous state-of-the-art 291 on ROSMAP and BRCA, applies a dynamic multimodal classification strategy. On Food-101 and 292 N24News, we compare against baseline unimodal networks (ViT and BERT/RoBERTa) and our 293 UniConcat baseline, where pre-trained image and text encoders are fine-tuned independently, and the unimodal representations are simply concatenated for classification. These are typical baselines used 295 in multimodal classification tasks. Detailed baseline descriptions are discussed in Appendix A.8. 296

4 Results

4.1 COMPARISON WITH BASELINES

The results are reported as the average and standard deviation over three runs on Food-101/N24News, and five runs on ROSMAP/BRCA. The best score is highlighted in bold, while the second-best score is underlined. The classification accuracy on N24News and Food-101 are displayed in Table 1 and 3 respectively. In the result tables, **ALI** denotes alignment (indicating if the method employs a contrastive component), while **AGG** specifies whether aggregation is early (combining unimodal feature) or late fusion (combining unimodal decisions).

307 The experimental results from Table 1, 2, 3, reveal the following findings: (i) M3CoL consistently 308 outperforms all SOTA methods across all text sources on N24News when using the same encoders, 309 beats SOTA on all evaluation metrics on ROSMAP and BRCA, and also achieves competitive results on Food-101; (ii) contrastive-based methods with any form of alignment demonstrate superior 310 performance compared to other multimodal methods; (iii) our proposed M3CoL method, which 311 employs a contrastive-based approach with shared alignment, improves over the traditional contrastive-312 based models and the latest SOTA multimodal methods. We visualize the unimodal and combined 313 representation distribution of our proposed method using UMAP plots in Figure 7 in Appendix A.6. 314

315 316

297

298 299

300

4.2 ANALYSIS OF OUR METHOD

Effect of Vanilla Mixup. Mixup involves two main components: the random convex combination of raw inputs and the corresponding convex combination of one-hot label encodings. To assess the performance of our M3CoL method in comparison to this Mixup strategy, we conduct experiments on Food-101 and N24News (text source: abstract). We remove the contrastive loss from our framework (Eq. 16) while keeping the rest of the modules unchanged. Table 4 shows that the Mixup technique underperforms relative to our proposed M3CoL approach (refer test accuracy plots illustrated in Figure 6a). The observed accuracy gap can be attributed to excessive noise introduced by label mixing, and the lack of a contrastive approach with an alignment component. This indicates that the

Method	Fus	ion	Bao	ckbone		ACC ↑	
	AGG	ALI	Image	Text	Headline	Caption	Abstract
Image-only	-	-	ViT	-	54.1 (n	o text source	e used)
Text-only	-	-	-	BERT	72.1	72.7	78.3
UniConcat	Early	X	ViT	BERT	78.6	76.8	80.8
UniS-MMC	Early	1	ViT	BERT	80.3	77.5	83.2
M3CoL (Ours)	Early	\checkmark	ViT	BERT	$80.8_{\pm_{0.05}}$	$78.0_{\pm_{0.03}}$	$83.8_{\pm_{0.06}}$
Text-only	-	-	-	RoBERTa	71.8	72.9	79.7
UniConcat	Early	X	ViT	RoBERTa	78.9	77.9	83.5
N24News	Early	X	ViT	RoBERTa	79.41	77.45	83.33
UniS-MMC	Early	1	ViT	RoBERTa	80.3	78.1	84.2
M3CoL (Ours)	Early	1	ViT	RoBERTa	$80.9_{\pm 0.19}$	$79.2_{\pm 0.08}$	$84.7_{\pm 0.03}$
	•				-0.10	=0.00	_0.00

Table 1: Accuracy (ACC) on N24News on three text sources. AGG denotes early/late modality fusion, ALI indicates presence/absence of alignment. Our method consistently outperforms SOTA across all text sources and backbone combinations. Baseline details are provided in Appendix A.8.

Method	Fus	ion		ROSMAP			BRCA		
	AGG	ALI	ACC ↑	F1 ↑	AUC ↑	ACC ↑	WF1 ↑	MF1 ↑	
GRidge	Early	X	76.0	76.9	84.1	74.5	72.6	65.6	
BPLSDA	Early	X	74.2	75.5	83.0	64.2	53.4	36.9	
BSPLSDA	Early	X	75.3	76.4	83.8	63.9	52.2	35.1	
MOGONET	Late	X	81.5	82.1	87.4	82.9	82.5	77.4	
TMC	Late	X	82.5	82.3	88.5	84.2	84.4	80.6	
CF	Early	X	78.4	78.8	88.0	81.5	81.5	77.1	
GMU	Early	X	77.6	78.4	86.9	80.0	79.8	74.6	
MOSEGCN	Early	X	83.0	82.7	83.2	86.7	86.8	81.1	
DYNAMICS	Early	X	85.7	86.3	<u>91.1</u>	<u>87.7</u>	<u>88.0</u>	<u>84.5</u>	
M3CoL (Ours)	Early	1	88.7 _{±0.94}	88.5 _{±0.94}	92.6 _{±0.59}	88.4 _{±0.57}	89.0 ±0.42	86.2 _{±0.54}	

Table 2: Comparison of Accuracy (ACC), Area Under the Curve (AUC), F1 score (F1) on ROSMAP, and Accuracy (ACC), Weighted F1 score (WF1), and Micro F1 score (MF1) on BRCA datasets. AGG denotes early/late modality fusion, ALI indicates presence/absence of alignment. Our method significantly outperforms SOTA across all metrics. Baseline details are provided in Appendix A.8.

> vanilla Mixup strategy introduces additional noise which impairs the model's ability to learn effective representations, while our M3CoL framework benefits from the structured contrastive approach.

Effect of Loss & Unimodality Supervision. To assess the necessity of each component in the framework, we investigate several design choices: (i) the framework's performance without the supervision of unimodal modules during training, and (ii) the performance differences between using only MultiSClip and only M3Co loss during end-to-end training. The M3CoL (No Unimodal Supervision) result indicates that excluding the unimodal prediction module results in a decline in performance as shown in Table 4 and Figure 6a, highlighting its importance as it allows the model to compensate for the weaknesses of one modality with the strengths of another. Additionally, the M3Co loss (only M3Co) outperforms the MultiSClip loss (only MultiSClip) by learning more robust representations through Mixup-based techniques, which prevent trivial discrimination of positive pairs. Furthermore, using an individual contrastive alignment approach (only M3Co) throughout the entire training process without transitioning to the MultiSClip loss results in suboptimal outcomes. This can be attributed to the risk of over-training with Mixup-based loss, which may negatively impact generalization. This demonstrates the necessity of the transition of the contrastive loss during training (0.33 M3Co + 0.67 MultiSClip). Figure 6b displays the accuracy plots on the N24News dataset, for these losses.

Visualization of Attention Heatmaps. The attention heatmaps generated using the embeddings from our trained M3CoL model in Figure 3 and 4 highlight image regions most relevant to the input

Method	Fus	ion	Backbo	ne	
memou	AGG	ALI	Image	Text	nee
Image-only	-	-	ViT	-	73.1
Text-only	-	-	-	BERT	86.8
UniConcat	Early	X	ViT	BERT	93.7
MCCE	Early	X	DenseNet	BERT	91.3
CentralNet	Early	X	LeNet5	LeNet5	91.5
GMU	Early	X	RNN	VGG	90.6
ELS-MMC	Early	X	ResNet-152	BOW features	90.8
MMBT	Early	X	ResNet-152	BERT	91.7
HUSE	Early	1	Graph-RISE	BERT	92.3
VisualBERT	X	1	FasterRCNN+BERT	BERT	92.3
PixelBERT	Early	1	ResNet	BERT	92.6
ViLT	Early	1	ViT	BERT	92.9
CMA-CLIP	Early	1	ViT	BERT	93.1
ME	Early	X	DenseNet	BERT	94.7
UniS-MMC	Early	\checkmark	ViT	BERT	94.7
M3CoL (Ours)	Early	1	ViT	BERT	$94.3_{\pm 0.04}$

Table 3: Accuracy (ACC) comparison on Food-101. AGG denotes early/late modality fusion, ALI indicates presence/absence of alignment. Baseline details are provided in Appendix A.8.

Method		AC	C↑	
	ROSMAP	BRCA	Food-101	N24News
Mixup	$84.13_{\pm_{0.74}}$	$84.52_{\pm 0.46}$	$93.14_{\pm_{0.02}}$	$81.57_{\pm_{0.24}}$
M3CoL (No Unimodal Supervision)	$85.14_{\pm 0.85}$	$86.93_{\pm 0.52}$	$94.12_{\pm 0.02}$	$84.26_{\pm 0.11}$
M3CoL (only MultiSClip)	$86.84_{\pm 0.34}$	$87.38_{\pm_{0.41}}$	$94.23_{\pm 0.01}$	$84.06_{\pm 0.18}$
M3CoL (only M3Co)	$87.42_{\pm 0.63}$	$87.74_{\pm 0.42}$	$94.24_{\pm 0.12}$	$84.57_{\pm 0.08}$
M3CoL (0.33 M3Co + 0.67 MultiSClip)	$88.67_{\pm_{0.94}}$	$\textbf{88.38}_{\pm_{0.57}}$	$\textbf{94.27}_{\pm_{0.04}}$	$84.72_{\pm_{0.03}}$

Table 4: Accuracy (ACC) on ROSMAP, BRCA, N24News, and Food-101 datasets under different settings of our method. For N24News, source: abstract and encoder: RoBERTa.

word. We generate text embeddings for class label words and corresponding image patch embeddings, computing attention scores as their dot product. This visualization aids in understanding the model's focus, decision-making process, and association between class labels and specific image regions. Importantly, it also indicates the correctness of the learned multimodal representations, revealing the model's ability to learn shared relations amongst different modalities, and ground visual concepts to semantically meaningful regions.



Figure 3: Text-guided visual grounding with varying input prompts. (a, e) Original images. (b-d)
Attention heatmaps for "ice cream" class. (f-h) Heatmaps for "falafel" class. Ice cream example: (b)
"Ice cream": Concentrated focus on ice cream, (c) "Cream": Maintained but diffused focus, (d) "Ice": Dispersed attention. Falafel example: (f) "Falafel": Localized focus on falafel, (g) "Salad": Attention shift to salad component, (h) "Rice": Minimal attention (absent in image). Warmer colors indicate higher attention scores.



Figure 4: Text-guided visual grounding with ablated model variations. (a) Original image. (b-f) Attention heatmaps generated using text embedding (class name: "Risotto") and patch embeddings for different variations of the model. Our proposed M3CoL model (f) demonstrates superior attention localization compared to ablated versions (b-e), corroborating the quantitative results presented in Table 4. Warmer colors indicate higher attention scores. (Here, No Unim: No Unimodal Supervision)

444 Testing on Random Data and Single-Corrupt Modal-445 ities. To showcase the benefits of our framework over 446 traditional contrastive methods, we evaluate the impact of 447 incorporating Mixup-based contrastive loss (M3Co) during training, highlighting its improvements over standard 448 approaches. It is well-established that deep networks tend 449 to exhibit overconfidence, particularly when making pre-450 dictions on random or adversarial inputs (Hendrycks & 451 Gimpel, 2016). Previous research has demonstrated that 452 Mixup can mitigate this issue, and our goal is to validate



Figure 5: N24News - Confidence scores when tested on random inputs.

its effectiveness in this context (Thulasidasan et al., 2019). We evaluate the confidence scores
 produced using M3CoL (0.33 M3Co + 0.67 MultiSClip) loss in comparison to only MultiSClip
 loss when predicting on random noise images and text encoder outputs. Our results show that the
 model trained with M3CoL exhibits lower confidence in its predictions when both modalities are
 replaced with random inputs. This demonstrates that incorporating M3CoL enhances the reliability
 of predictions, especially in the presence of corrupted or random inputs.

459 To evaluate the robustness of our approach, we conduct 460 experiments where one input modality was corrupted 461 with random noise. Table 5 compares the performance 462 of M3CoL (0.33 M3Co + 0.67 MultiSClip) against only 463 MultiSClip under these conditions. Our M3CoL method 464 demonstrates superior robustness to modality corruption, consistently outperforming MultiSClip. For image corrup-465 tion, we substituted the original images with random noise 466

Method	Modality Corrupted				
	Image	Text			
MutliSClip	76.06	46.91			
M3CoL	77.24	47.94			

Table 5: N24News - Accuracy when tested on data with one corrupt modality.

sampled from a Gaussian distribution, parameterized to match the mean and variance of the training
set. Similarly, for text corruption, we replaced the original text embeddings with random outputs
from the text encoder, again following a Gaussian distribution with statistics matching the training
data. For both the above experiments, we use the N24News dataset, with the abstract as the text
source and a RoBERTa-based text encoder.

Error Analysis. To evaluate the efficacy of our multimodal approach in integrating and leveraging 472 image and text features, we performed a comprehensive error analysis, comparing it with image-only 473 (ViT) and text-only (RoBERTa) models using the N24News dataset (refer Table 9 in Appendix A.5). 474 The analysis reveals that our method excels when both modalities are correctly classified (42.71:0.03 475 correct-to-incorrect ratio). This demonstrates that our model can learn valuable insights from the 476 fusion of image and text features, which may not be discovered when processing them separately. In 477 cases where only one modality is correctly classified, our model effectively leverages the accurate 478 modality (27.77+8.11=35.88):(1.29+3.25=4.54) correct-to-incorrect ratio. This demonstrates our 479 method's robustness and its ability to outperform unimodal approaches.

480 481

437

438

439

440

441

442 443

5 RELATED WORK

482 483

484 Contrastive Learning. Contrastive learning has driven significant progress in unimodal and multi 485 modal representation learning by distinguishing between similar (positive) and dissimilar (negative) pairs. In multimodal contexts, cross-modal contrastive techniques align representations from different

486 modalities (Radford et al., 2021; Jia et al., 2021; Kamath et al., 2021), with approaches like CrossCLR 487 (Zolfaghari et al., 2021) and GMC (Poklukar et al., 2022) focusing on global and modality-specific 488 representations. Contrastive learning approaches for paired image-text data, such as CLIP (Radford 489 et al., 2021), ALIGN (Jia et al., 2021), and BASIC (Pham et al., 2023), have demonstrated remarkable 490 success across diverse vision-language tasks. Subsequent works have aimed to enhance the efficacy and data efficiency of CLIP training, incorporating self-supervised techniques (SLIP (Mu et al., 2022), 491 DeCLIP (Li et al., 2021)) and fine-grained alignment (FILIP (Li et al., 2023)). The CLIP framework 492 relies on data augmentations to prevent overfitting and the learning of ineffective shortcuts (Chen 493 et al., 2020; Robinson et al., 2021), a common practice in contrastive learning. 494

495 Unimodal and Multimodal Data Augmentation. Data augmentation has been integral to the success of deep learning, especially for small training sets. In computer vision, techniques have 496 evolved from basic transformations to advanced methods like Cutout (DeVries & Taylor, 2017), 497 Mixup (Zhang et al., 2017), CutMix (Yun et al., 2019), and automated approaches (Cubuk et al., 498 2019; 2020). NLP augmentation includes paraphrasing, token replacement (Zhang et al., 2015; Jiao 499 et al., 2019), and noise injection (Yan et al., 2019). Multimodal data augmentation, primarily focused 500 on vision-text tasks, has seen limited exploration, with approaches including back-translation for 501 visual question answering (Tang et al., 2020), text generation from images (Wang et al., 2021b), and 502 external knowledge querying for cross-modal retrieval (Gur et al., 2021). MixGen (Hao et al., 2023) 503 generates new image-text pairs through image interpolation and text concatenation. In contrast, our 504 proposed augmentation technique focusing on the early training phase is fully automatic, applicable 505 to arbitrary modalities, and designed to leverage inherent shared relations in multimodal data. 506

Relation to Mixup. Mixup (Zhang et al., 2017), a pivotal regularization strategy, enhances model robustness and generalization by generating synthetic samples through convex combinations of existing data points. Originally introduced for computer vision, it has been adapted to NLP by applying the technique to text embeddings (Guo et al., 2019). Our proposed augmentation differs from Mixup in several key aspects: it is designed for multi-modal data, takes inputs from different modalities, and does not rely on one-hot label encodings. By extending the Mixup paradigm to complex, multi-modal scenarios and focusing on the early training phase, our method broadens its applicability while leveraging inherent shared relations in multimodal data.

- 514 515
- 516
- 517

6 CONCLUSION

- 518 519
- 520

Aligning representations across modalities presents significant challenges due to the complex, often 521 non-bijective relationships in real-world multimodal data (Liang et al., 2022a). These relationships 522 can involve many-to-many mappings or even lack clear associations, as exemplified by linguistic 523 ambiguities and synonymy in vision-language tasks. We propose M3Co, a novel contrastive-based 524 alignment method that captures shared relations beyond explicit pairwise associations by aligning 525 mixed samples from one modality with corresponding samples from others. Our approach incor-526 porates Mixup-based contrastive learning, introducing controlled noise that mirrors the inherent 527 variability in multimodal data, thus enhancing robustness and generalizability. The M3Co loss, 528 combined with an architecture leveraging unimodal and fusion modules, enables continuous updating 529 of representations necessary for accurate predictions and deeper integration of modalities. Our 530 method generalizes across diverse domains, including image-text, high-dimensional multi-omics, and data with more than two modalities. Experiments on four public multimodal classification 531 datasets demonstrate the effectiveness of our approach in learning robust representations that surpass 532 traditional multimodal alignment techniques. 533

M3CoL demonstrates promising results, yet faces optimization challenges due to the inherent
limitations of multimodal frameworks, particularly extended training times on large-scale datasets
like Food-101. The method's modality-agnostic nature and effective use of Mixup augmentation
suggest its potential adaptability to various multimodal tasks, especially where data augmentation
and learning real-world inter-modal relationships are crucial. Our method can easily be adapted to
multimodal tasks such as visual question answering and information retrieval. These advancements
will likely broaden M3CoL's impact in multimodal research.

540 REFERENCES 541

547

551

576

542	Mahdi Abavisani, Liwei Wu, Shengli Hu, Joel Tetreault, and Alejandro Jaimes. Multimodal catego-
543	rization of crisis events in social media. In Proceedings of the IEEE/CVF Conference on Computer
544	Vision and Pattern Recognition, pp. 14679–14689, 2020.

John Arevalo, Thamar Solorio, Manuel Montes-y Gómez, and Fabio A González. Gated multimodal 546 units for information fusion. arXiv preprint arXiv:1702.01992, 2017.

Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: 548 A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2): 549 423-443, 2018. 550

- Olivier Chapelle, Jason Weston, Léon Bottou, and Vladimir Vapnik. Vicinal risk minimization. 552 Advances in neural information processing systems, 13, 2000.
- 553 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for 554 contrastive learning of visual representations. In International conference on machine learning, pp. 555 1597-1607. PMLR, 2020. 556
- Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: 558 Learning augmentation strategies from data. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 113–123, 2019. 559
- Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated 561 data augmentation with a reduced search space. In Proceedings of the IEEE/CVF conference on 562 computer vision and pattern recognition workshops, pp. 702–703, 2020. 563
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018. 565
- 566 Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks 567 with cutout. arXiv preprint arXiv:1708.04552, 2017. 568
- 569 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An 570 image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint 571 arXiv:2010.11929, 2020. 572
- 573 Yuting Gao, Jinfeng Liu, Zihan Xu, Jun Zhang, Ke Li, Rongrong Ji, and Chunhua Shen. Pyramid-574 clip: Hierarchical feature alignment for vision-language model pretraining. Advances in neural 575 information processing systems, 35:35959–35970, 2022.
- Yuting Gao, Jinfeng Liu, Zihan Xu, Tong Wu, Enwei Zhang, Ke Li, Jie Yang, Wei Liu, and Xing 577 Sun. Softclip: Softer cross-modal alignment makes clip stronger. In Proceedings of the AAAI 578 Conference on Artificial Intelligence, volume 38, pp. 1860–1868, 2024. 579
- 580 Aditya Sharad Golatkar, Alessandro Achille, and Stefano Soatto. Time matters in regularizing deep networks: Weight decay and data augmentation affect early learning dynamics, matter little near 581 convergence. Advances in Neural Information Processing Systems, 32, 2019. 582
- 583 Hongyu Guo, Yongyi Mao, and Richong Zhang. Augmenting data with mixup for sentence classifica-584 tion: An empirical study. arXiv preprint arXiv:1905.08941, 2019. 585
- Shir Gur, Natalia Neverova, Chris Stauffer, Ser-Nam Lim, Douwe Kiela, and Austin Reiter. Cross-586 modal retrieval augmentation for multi-modal classification. arXiv preprint arXiv:2104.08108, 587 2021. 588
- Zongbo Han, Changqing Zhang, Huazhu Fu, and Joey Tianyi Zhou. Trusted multi-view classification. 590 In International Conference on Learning Representations, 2020. 591
- Zongbo Han, Fan Yang, Junzhou Huang, Changqing Zhang, and Jianhua Yao. Multimodal dynamics: 592 Dynamical fusion for trustworthy multimodal classification. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 20707-20717, 2022.

594 595 596	Xiaoshuai Hao, Yi Zhu, Srikar Appalaraju, Aston Zhang, Wanqian Zhang, Bo Li, and Mu Li. Mixgen: A new multi-modal data augmentation. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pp. 379–389, 2023.
597 598 599	Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. <i>arXiv preprint arXiv:1610.02136</i> , 2016.
600 601 602	Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshmi- narayanan. Augmix: A simple data processing method to improve robustness and uncertainty. <i>arXiv preprint arXiv:1912.02781</i> , 2019.
603 604 605	Danfeng Hong, Lianru Gao, Naoto Yokoya, Jing Yao, Jocelyn Chanussot, Qian Du, and Bing Zhang. More diverse means better: Multimodal deep learning meets remote-sensing imagery classification. <i>IEEE Transactions on Geoscience and Remote Sensing</i> , 59(5):4340–4354, 2020.
607 608 609	Yu Huang, Chenzhuang Du, Zihui Xue, Xuanyao Chen, Hang Zhao, and Longbo Huang. What makes multi-modal learning better than single (provably). <i>Advances in Neural Information Processing Systems</i> , 34:10944–10956, 2021.
610 611	Zhicheng Huang, Zhaoyang Zeng, Bei Liu, Dongmei Fu, and Jianlong Fu. Pixel-bert: Aligning image pixels with text by deep multi-modal transformers. <i>arXiv preprint arXiv:2004.00849</i> , 2020.
612 613 614 615 616	Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In <i>International conference on machine learning</i> , pp. 4904–4916. PMLR, 2021.
617 618 619	Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Tinybert: Distilling bert for natural language understanding. <i>arXiv preprint arXiv:1909.10351</i> , 2019.
620 621 622	Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In <i>Proceedings of the</i> <i>IEEE/CVF International Conference on Computer Vision</i> , pp. 1780–1790, 2021.
623 624 625	Douwe Kiela, Edouard Grave, Armand Joulin, and Tomas Mikolov. Efficient large-scale multi-modal classification. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 32, 2018.
626 627	Douwe Kiela, Suvrat Bhooshan, Hamed Firooz, Ethan Perez, and Davide Testuggine. Supervised multimodal bitransformers for classifying images and text. <i>arXiv preprint arXiv:1909.02950</i> , 2019.
628 629 630	Sungnyun Kim, Gihun Lee, Sangmin Bae, and Se-Young Yun. Mixco: Mix-up contrastive learning for visual representation. <i>arXiv preprint arXiv:2010.06300</i> , 2020.
631 632 633	Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convo- lution or region supervision. In <i>International conference on machine learning</i> , pp. 5583–5594. PMLR, 2021.
634 635	Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. <i>arXiv preprint arXiv:1412.6980</i> , 2014.
637 638	Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. <i>arXiv preprint arXiv:1908.03557</i> , 2019.
639 640 641	Yangguang Li, Feng Liang, Lichen Zhao, Yufeng Cui, Wanli Ouyang, Jing Shao, Fengwei Yu, and Junjie Yan. Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm. <i>arXiv preprint arXiv:2110.05208</i> , 2021.
642 643 644 645	Yanghao Li, Haoqi Fan, Ronghang Hu, Christoph Feichtenhofer, and Kaiming He. Scaling language- image pre-training via masking. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision</i> <i>and Pattern Recognition</i> , pp. 23390–23400, 2023.
646 647	Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. Foundations and trends in multimodal machine learning: Principles, challenges, and open questions. <i>arXiv preprint arXiv:2209.03430</i> , 2022a.

648 649 650 651	Tao Liang, Guosheng Lin, Mingyang Wan, Tianrui Li, Guojun Ma, and Fengmao Lv. Expanding large pre-trained unimodal models with multimodal information injection for image-text multimodal classification. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 15492–15501, 2022b.
653 654 655	Huidong Liu, Shaoyuan Xu, Jinmiao Fu, Yang Liu, Ning Xie, Chien-Chih Wang, Bryan Wang, and Yi Sun. Cma-clip: Cross-modality attention clip for image-text classification. <i>arXiv preprint arXiv:2112.03562</i> , 2021.
656 657 658	Zixuan Liu, Ziqiao Wang, Hongyu Guo, and Yongyi Mao. Over-training with mixup may hurt generalization. <i>arXiv preprint arXiv:2303.01475</i> , 2023.
659 660 661	Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. Slip: Self-supervision meets language-image pre-training. In <i>European conference on computer vision</i> , pp. 529–544. Springer, 2022.
662 663 664	Pradyumna Narayana, Aniket Pednekar, Abishek Krishnamoorthy, Kazoo Sone, and Sugato Basu. Huse: Hierarchical universal semantic embeddings. <i>arXiv preprint arXiv:1911.05978</i> , 2019.
665 666	Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. <i>arXiv preprint arXiv:1807.03748</i> , 2018.
667 668 669 670	Hieu Pham, Zihang Dai, Golnaz Ghiasi, Kenji Kawaguchi, Hanxiao Liu, Adams Wei Yu, Jiahui Yu, Yi-Ting Chen, Minh-Thang Luong, Yonghui Wu, et al. Combined scaling for zero-shot transfer learning. <i>Neurocomputing</i> , 555:126658, 2023.
671 672 673	Petra Poklukar, Miguel Vasco, Hang Yin, Francisco S Melo, Ana Paiva, and Danica Kragic. Geometric multimodal contrastive representation learning. In <i>International Conference on Machine Learning</i> , pp. 17782–17800. PMLR, 2022.
674 675 676 677	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 <i>International conference on machine learning</i> , pp. 8748–8763. PMLR, 2021.
679 680 681	Joshua Robinson, Li Sun, Ke Yu, Kayhan Batmanghelich, Stefanie Jegelka, and Suvrit Sra. Can contrastive learning avoid shortcut solutions? <i>Advances in neural information processing systems</i> , 34:4974–4986, 2021.
682 683 684	Zhiqiang Shen, Zechun Liu, Zhuang Liu, Marios Savvides, Trevor Darrell, and Eric Xing. Un-mix: Rethinking image mixtures for unsupervised visual representation learning. In <i>Proceedings of the</i> <i>AAAI Conference on Artificial Intelligence</i> , volume 36, pp. 2216–2224, 2022.
686 687 688	Amrit Singh, Casey P Shannon, Benoît Gautier, Florian Rohart, Michaël Vacher, Scott J Tebbutt, and Kim-Anh Lê Cao. Diablo: an integrative approach for identifying key molecular drivers from multi-omics assays. <i>Bioinformatics</i> , 35(17):3055–3062, 2019.
689 690 691	Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. Advances in neural information processing systems, 29, 2016.
692 693 694 695	Ruixue Tang, Chao Ma, Wei Emma Zhang, Qi Wu, and Xiaokang Yang. Semantic equivalent adversarial data augmentation for visual question answering. In <i>Computer Vision–ECCV 2020:</i> 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIX 16, pp. 437–453. Springer, 2020.
696 697 698 699	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> , 2023.
700 701	Sunil Thulasidasan, Gopinath Chennupati, Jeff A Bilmes, Tanmoy Bhattacharya, and Sarah Michalak. On mixup training: Improved calibration and predictive uncertainty for deep neural networks. <i>Advances in neural information processing systems</i> , 32, 2019.

702 703 704	Mark A Van De Wiel, Tonje G Lien, Wina Verlaat, Wessel N van Wieringen, and Saskia M Wilting. Better prediction by use of co-data: adaptive group-regularized ridge regression. <i>Statistics in medicine</i> , 35(3):368–381, 2016.
705 706	Vikas Verma, Alex Lamb, Christopher Beckham, Amir Najafi, Ioannis Mitliagkas, David Lopez-Paz,
707	and Yoshua Bengio. Manifold mixup: Better representations by interpolating hidden states. In
708	International conference on machine learning, pp. 6438–6447. PMLR, 2019.
709	Valantin Vialzauf, Alavis Lacheruy, Stánhana Patauy, and Frádáric Juria, Cantralnat, a multilayar
710	approach for multimodal fusion. In Proceedings of the European Conference on Computer Vision
711	(ECCV) Workshops, pp. 0–0, 2018.
712	
713	Jiahui Wang, Nanqing Liao, Xiaofei Du, Qingfeng Chen, and Bizhong Wei. A semi-supervised
714 715	approach for the integration of multi-omics data based on transformer multi-head self-attention mechanism and graph convolutional networks. <i>BMC genomics</i> , 25(1):86, 2024.
716	Tongyin Wang, Wei Shao, Zhi Huang, Haiyu Tang, Jie Zhang, Zhengming Ding, and Kun Huang
717 718	Mogonet integrates multi-omics data using graph convolutional networks allowing patient classifi- cation and biomarker identification. <i>Nature communications</i> , 12(1):3445, 2021a.
719	
720	Xin Wang, Devinder Kumar, Nicolas Thome, Matthieu Cord, and Frederic Precioso. Recipe recogni-
721 722	Expo Workshops (ICMEW), pp. 1–6. IEEE, 2015.
723	Zhen Wang, Xu Shan, Xiangxie Zhang, and Jie Yang, N24news: A new dataset for multimodal
724	news classification. In Proceedings of the Language Resources and Evaluation Conference, pp.
725	6768-6775, Marseille, France, June 2022. European Language Resources Association. URL
726	https://aclanthology.org/2022.lrec-1.729.
727	Zixu Wang, Yishu Miao, and Lucia Specia. Cross-modal generative augmentation for visual question
728	answering. arXiv preprint arXiv:2105.04780, 2021b.
730	Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-
731 732	parametric instance discrimination. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 3733–3742, 2018.
733	Ge Yan, Yu Li, Shu Zhang, and Zhenyu Chen. Data augmentation for deep learning of judgment
735 736 737	documents. In Intelligence Science and Big Data Engineering. Big Data and Machine Learning: 9th International Conference, IScIDE 2019, Nanjing, China, October 17–20, 2019, Proceedings, Part II 9, pp. 232–242. Springer, 2019.
738 739 740	Hao Yu, Huanyu Wang, and Jianxin Wu. Mixup without hesitation. In <i>Image and Graphics: 11th International Conference, ICIG 2021, Haikou, China, August 6–8, 2021, Proceedings, Part II 11</i> , pp. 143–154. Springer, 2021.
741 742 743 744	Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In <i>Proceedings</i> of the IEEE/CVF international conference on computer vision, pp. 6023–6032, 2019.
745 746 747	Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. <i>arXiv preprint arXiv:1710.09412</i> , 2017.
748 749	Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. <i>Advances in neural information processing systems</i> , 28, 2015.
750 751 752 753	Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis P Langlotz. Con- trastive learning of medical visual representations from paired images and text. In <i>Machine</i> <i>Learning for Healthcare Conference</i> , pp. 2–25. PMLR, 2022.
754 755	Liu Zhuang, Lin Wayne, Shi Ya, and Zhao Jun. A robustly optimized bert pre-training approach with post-training. In <i>Proceedings of the 20th chinese national conference on computational linguistics</i> , pp. 1218–1227, 2021.

Mohammadreza Zolfaghari, Yi Zhu, Peter Gehler, and Thomas Brox. Crossclr: Cross-modal contrastive learning for multi-modal video representations. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 1450–1459, 2021. Heqing Zou, Meng Shen, Chen Chen, Yuchen Hu, Deepu Rajan, and Eng Siong Chng. Unis-mmc: Multimodal classification via unimodality-supervised multimodal contrastive learning. arXiv preprint arXiv:2305.09299, 2023.

810 A APPENDIX

812 A.1 EXPERIMENTAL DETAILS

The models were trained on either an NVIDIA RTX A6000 or an NVIDIA A100-SXM4-80GB GPU. The results are reported as the average and standard deviation over three runs on Food-101 and N24News, and five runs on ROSMAP and BRCA. We use a grid search on the validation set to search for optimal hyperparameters. The temperature parameter for the M3Co and MultiSClip losses is set to 0.1. The corresponding loss coefficient β is 0.1 to keep the loss value in the same range as the other losses. We use the Adam optimizer (Kingma & Ba, 2014) for all datasets. For Food-101 and N24News, the learning rate scheduler is ReduceLROnPlateau with validation accuracy as the monitored metric, lr factor of 0.2, and lr patience of 2. For ROSMAP and BRCA, we use the StepLR scheduler with a step size of 250. For Food-101 and N24News, the maximum token length of the text input for the BERT/RoBERTa encoders is 512. Other hyperparameter details are provided in Table 6.

Hyperparameter	N24News	Food-101	ROSMAP	BRCA
Embedding dimension	768	768	1000	768
Classifier dimension	256	256	1000	768
Learning rate	10^{-4}	10^{-4}	$5 \cdot 10^{-3}$	$5\cdot 10^{-3}$
Weight decay	10^{-4}	10^{-4}	10^{-3}	10^{-3}
Batch size	32	32	-	-
Batch gradient	128	128	-	-
Dropout (classifier)	0	0	0.5	0.5
Epochs	50	50	500	500

Table 6: Experimental hyperparameter values for our proposed model across all the four datasets.

Code. The code is attached as a zip file. Upon acceptance, the code and checkpoints will be made publicly available on GitHub.

A.2 DATASET INFORMATION AND SPLITS

The datasets used in our experiments can be downloaded from the following sources: Food 101 from https://visiir.isir.upmc.fr, N24News from https://github.com/
 billywzh717/N24News, and BRCA and ROSMAP from https://github.com/txWang/
 MOGONET.

To ensure a fair comparison with previous works, we adopt the default split method detailed in Table 7. As the Food-101 dataset does not include a validation set, we partition 5,000 samples from the training set to create one, which is consistent with other baselines.

Dataset	Modalities	Modality Types	Train	Validation	Test	Classes
Food-101	2	Image, text	60101	5000	21695	101
N24News	2	Image, text	48988	6123	6124	24
ROSMAP	3	mRNA, miRNA, DNA	245	-	106	2
BRCA	3	mRNA, miRNA, DNA	612	-	263	5

Table 7: Statistics for the four datasets: Food-101, N24News, ROSMAP, and BRCA. Note: miRNA stands for microRNA, and mRNA stands for messenger RNA.

A.3 ANALYSIS UNDER DIFFERENT MODEL VARIATIONS

In addition to the ACC scores presented in Table 4, we also report the performance of other metrics,
 where available, on the ROSMAP and BRCA datasets under various settings of our method, as shown
 in Table 8. This thorough evaluation supports our conclusion that each component of our framework
 is crucial for achieving optimal overall performance.

Method	ROS	MAP	BRCA		
Method	F1 ↑	AUC ↑	WF1 ↑	MF1 ↑	
Mixup	84.45+0.48	88.73+0.52	84.66+0.48	82.88+0.3	
M3CoL (No Unimodal Supervision)	$86.27_{\pm 0.48}$	$90.20_{\pm 0.83}$	$86.92_{\pm 0.36}$	$85.08_{\pm 0.3}$	
M3CoL (only MultiSClip)	$86.76_{\pm 0.58}$	$90.75_{\pm 0.40}$	$87.41_{\pm 0.47}$	$85.48_{\pm 0.3}$	
M3CoL (only M3Co)	$87.54_{\pm 0.62}$	$91.41_{\pm 0.54}$	$88.06_{\pm 0.41}$	$85.82_{\pm 0.3}$	
M3CoL (0.33 M3Co + 0.67 MultiSClip)	88.51 ±0.94	92.62 ±0.59	89.02 $_{\pm_{0.42}}$	86.20 _{±0.5}	

Table 8: Comparison of F1 score (F1), Area Under the Curve (AUC) on ROSMAP, and Weighted F1 score (WF1), Micro F1 score (MF1) on BRCA, under different settings of our method.

A.4 ABLATION STUDIES ON THE N24NEWS DATASET

878 The accuracy plots for the N24News dataset (text source: abstract, text encoder: RoBERTa) are used to compare our method and its variants. Our proposed M3CoL approach outperforms the 880 Mixup technique, as shown in Figure 6a. Ablating the unimodal supervision in M3CoL leads to a performance decline, indicating the importance of the unimodal prediction module, as shown in 882 Figure 6a. Furthermore, the M3Co loss achieves better results than the MultiSClip loss. Training solely with either the M3Co loss or the MultiSClip loss alignment approach yields suboptimal performance when compared to their strategic combination, as shown in Figure 6b. The quantitative results are given in Table 4.



(a) Comparison of M3CoL and its variants using Mixup and No Unimodal Supervision. 899

(b) Comparison of M3CoL and its variants using only M3Co and only MultiSoftClip loss.

Figure 6: Test accuracy plots showing comparison of M3CoL and its variants on the N24News dataset (text source: abstract, text encoder: RoBERTa).

902 903 904

905

906

907

908

909

910

873

874

875 876

877

879

881

883

884

885

890

891

892 893

896

897

900

901

A.5 ERROR ANALYSIS

We provide an in-depth error analysis on the N24News dataset in Section 4.2. As shown in Table 9, our method excels not only when both modalities are correctly classified but also demonstrates strong performance even when both are misclassified, indicating effective feature fusion. Additionally, it successfully leverages information when only one modality is classified correctly. This analysis demonstrates our method's robustness and highlights its superiority over unimodal approaches.

911 A.6 UMAP PLOTS 912

913 We generate UMAP plots on embeddings derived from the N24News and Food-101 datasets to 914 visualize the clustering performance of our M3CoL model. For each dataset, we randomly select 915 10 classes and generate the corresponding embeddings from the image encoder, text encoder, and their concatenated multimodal representatio, using our trained M3CoL model. Figure 7 shows that 916 the image embeddings depict less distinct clusters, indicating less effective inter-cluster separation 917 compared to the text encoder embeddings. The concatenated embeddings, however, result in the

Unimodal Prediction		Multimodal Prediction $\%$	
Text	Image	Correct	Incorrect
True	True	42.71	0.03
Гrue	False	27.77	1.29
False	True	8.11	3.25
False	False	2.31	14.53

Table 9: Error analysis on N24News with text encoder RoBERTa and text source "headline". True and False denote the correctness of the unimodal predictions. Multimodal Prediction % shows the resulting test set ratio of the final predictions.

best-defined clusters, suggesting that the final multimodal representations preserve and potentially enhance class-distinguishing features. These observations align with our quantitative results presented in Table 1 and 3.



Figure 7: UMAP plots of embeddings from the N24News (source: abstract and encoder: RoBERTa) and Food-101 datasets. We generate UMAP plots for the representations generated by the image encoder, text encoder, and their concatenated multimodal representations, using our trained M3CoL model. Concatenated embeddings exhibit superior clustering, while text embeddings outperform image embeddings. Consistent patterns across datasets demonstrate M3CoL's effectiveness in fusing multimodal information and enhancing semantic representations.

A.7 ADDITIONAL VISUALIZATION ATTENTION HEATMAPS

Following Section 4.2, we generate heatmaps using class embeddings and patch embeddings for some more examples in Food-101. These are depicted in Figure 8.

A.8 BASELINE DETAILS

The baselines used in our comparisons are described in details as follows:

- **GRidge** (Van De Wiel et al., 2016) dynamically incorporates multimodal data to adjust regularization penalties, improving predictive accuracy in genomic classification scenarios.
- **BPLSDA** (Block partial least squares discriminant analysis) (Singh et al., 2019) analyzes multimodal data in latent space and **BSPLSDA** (Block sparse partial least squares discrimi-



1000		
1026		and the ground truth, and using this insight to align feature vectors across various modalities
1027		through a contrastive loss.
1028		• ME (Liang et al., 2022b) leverages cross-modal information by transforming features
1029		between modalities. It achieves this by integrating a Multimodal Information Injection
1030		Plug-in (MI2P) with pre-trained models, enabling them to process image-text pairs without
1031		structural modifications.
1032		• MMPT (Viale at al. 2010) lower gras the strengths of prostrained text and image anecders
1033		• MINIDI (Kiela et al., 2019) leverages the strengths of pre-trained text and image encoders,
1034		into the textual taken space
1035		into the textual token space.
1036		• ELS-MMC (Kiela et al., 2018) investigates various multimodal fusion techniques for
1037		integrating discrete (text) and continuous (visual) modalities to enhance classification tasks
1038		in a resource-efficient manner.
1039		• N24News (Wang et al., 2022) presents a novel dataset from The New York Times with text
1040		and image data across 24 categories. It utilizes a multitask multimodal strategy, employing
10/11		ViT for image processing and RoBERTa for text analysis, with features concatenated for
1040		final classification.
1042		• CentralNet (Vielzeuf et al., 2018) uses separate convolutional networks for each modality.
1043		linked via a central network that generates a unified feature representation and also applies
1044		multi-task learning to refine and regulate these features.
1045		• CMU (Arayalo et al. 2017) employs multiplicative gates that dynamically adjust the
1046		influence of each modality on its activation thereby deriving a sophisticated intermediate
1047		representation tailored for specific applications
1048		
1049		• VISUAIBERT (Li et al., 2019) employs a series of transformer layers to align textual elements
1050		and corresponding image regions through self-attention mechanisms.
1051		• PixelBert (Huang et al., 2020) directly aligns image pixels with textual descriptions using a
1052		deep multi-modal transformer and establishes a direct semantic connection at the pixel and
1053		text level.
1054		• ViLT (Kim et al., 2021) implements a BERT-like transformer model that processes visual
1055		data in a convolution-free manner, similar to textual data, thereby simplifying input feature
1056		extraction and reducing computational demands.
1057		• HUSE (Narayana et al., 2019) constructs a shared latent space that aligns image and text
1058		embeddings based on their semantic similarity, enhancing cross-modal representation.
1059		• CMA CLID (Liu at al. 2021) anhances CLID (Badford at al. 2021) by integrating two
1060		• CMA-CLIF (Liu et al., 2021) eminances CLIF (Radiolu et al., 2021) by integrating two cross modality attention mechanisms: sequence wise and modality wise attention. These
1061		attention modules refine the relationships between image patches and text tokens, allowing
1062		the model to focus on relevant modalities for specific tasks
1063		MOSECCN (Ware at al. 2024) at il-
1064		• INIDEGUN (wang et al., 2024) utilizes transformer multi-head self-attention and Simi-
1065		information is then fed into a self ensembling Graph Convolutional Network (SECCN) for
1066		semi-supervised training and testing
1067		sonn-supervised daming and testing.
1069		
1060	A.9	USE OF GENERATIVE AI MODELS
1009	In thi	is work we use the following generative AI model:
1070	in ull	is work, we use the following generative fit model.
1071		• Gemini 1.0 Pro (Team et al., 2023) to generate food item images and captions as displayed
1072		in Figure 1, which serve as sample representations from the Food-101 dataset.
1073		
1074		
1075		
1076		
1077		
1078		
1079		