# Understanding and Constructing Latent Modality Structures in Multi-Modal Representation Learning

Qian Jiang[1], Changyou Chen[2,3], Han Zhao[1,3], Liqun Chen[*], Qing Ping[3],
Son Dinh Tran[3], Yi Xu[3], Belinda Zeng[3], Trishul Chilimbi[3]
[1]University of Illinois at Urbana-Champaign    [2]University at Buffalo    [3]Amazon

qianj3@illinois.edu    lqchen06@outlook.com

{vchencha, uhanzhao, pingqing, sontran, yxaamzn, zengb, trishulc}@amazon.com

## Abstract

*Contrastive loss has been increasingly used in learning representations from multiple modalities. In the limit, the nature of the contrastive loss encourages modalities to exactly match each other in the latent space. Yet it remains an open question how the modality alignment affects the downstream task performance. In this paper, based on an information-theoretic argument, we first prove that exact modality alignment is sub-optimal in general for downstream prediction tasks. Hence we advocate that the key of better performance lies in meaningful latent modality structures instead of perfect modality alignment. To this end, we propose three general approaches to construct latent modality structures. Specifically, we design 1) a deep feature separation loss for intra-modality regularization; 2) a Brownian-bridge loss for inter-modality regularization; and 3) a geometric consistency loss for both intra- and inter-modality regularization. Extensive experiments are conducted on two popular multi-modal representation learning frameworks: the CLIP-based two-tower model and the ALBEF-based fusion model. We test our model on a variety of tasks including zero/few-shot image classification, image-text retrieval, visual question answering, visual reasoning, and visual entailment. Our method achieves consistent improvements over existing methods, demonstrating the effectiveness and generalizability of our proposed approach on latent modality structure regularization.*

## 1. Introduction

Vision-language representation learning aims to learn generic representations from images and texts that could benefit multimodal downstream applications. As the two modalities are essentially from different data sources and distributions, how to effectively fuse the two modalities has
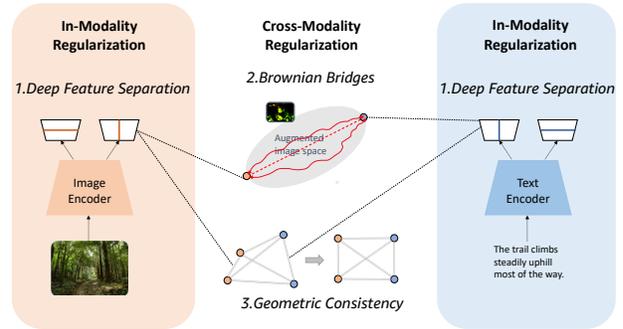
---

*Work done while at Amazon.



Figure 1. Constructing latent modality structures to improve multi-modal representation learning.

become an important question. Some work aims to unify the representations of two modalities in one encoder, where the image and text are usually tokenized into sequences [58, 59, 63, 64]. Another line of research represents the image and text modality separately with modality-specific encoders and utilizes contrastive learning to align the modalities, achieving state-of-the-art performance on multiple downstream applications [13, 25, 30, 31, 40, 48, 52, 53, 67].

Despite the successful empirical practice of contrastive loss in multi-modal representation learning, it remains an open question whether bridging and aligning the two modalities always brings benefits to downstream tasks. One concept closely related to this question is the modality gap [34, 48, 66, 69], where it is defined as the distance between the feature distributions of the two modalities. Modality alignment can be considered as reducing the modality gap. At a first glance, one would conjecture that contrastive loss would reduce the modality gap by pulling positive (paired) image and text data together for better representation. However, a recent study [34] shows evidence that contrastive learning does not always reduce the modality gap. Furthermore, we also show in our empirical analysis that a reduced modality gap does not always guaran-

tee better performance in downstream applications. Motivated by these empirical observations, in this paper we first theoretically study the modality gap problem, by showing that when the modality gap is zero, i.e., exact alignment between the two modalities, the learned representations necessarily have to pay a price for the downstream prediction task, which we term as the *information gap* between the two modalities (Theorem 3.1). Intuitively, this is because that representations with zero modality gap can only preserve predictive information present in *both* of the modalities at the cost of losing the *modality-specific* information.

Our theory then suggests that instead of exact modality matching, whether learned representations are meaningful is an important factor in multi-modal representation learning. In particular, we propose to improve on top of contrastive learning with regularizations to construct better latent structures. We consider intra-modality, inter-modality, and intra-inter-modality regularizations. These regularizations are generalizable and can be applied to various vision-language models with modality-specific encoders. Specifically, for intra-modality regularization, motivated by our theoretic result, we propose deep feature separation to encourage the model to preserve both the modality-shared and modality-specific information in different components. For inter-modality regularization, we aim to bridge two modalities with their augmentations. Consequently, we proposed a Brownian bridge loss between the triplet of (text, augmented image, image) to regularize the inter-modality structures. For intra-inter-modality regularization, we introduce the geometric consistency loss that promotes geometric symmetry in the latent space. In summary, the main contributions of this paper are:

- We conduct empirical and theoretical analysis on understanding the impact of the modality alignment on downstream tasks. We show that a reduced modality gap does not always guarantee better performance, and can instead hurt the performance when the information gap between the two modalities is large (Theorem 3.1). Combined with the existing theory of contrastive learning, our theory suggests preserving both modality-shared and modality-specific information.

- Inspired by our theory, we propose three instrumental regularizations on top of the contrastive loss, *i.e.*, the intra-modality, inter-modality, and intra-inter-modality regularizations to improve latent modality structures.

- We conduct extensive and comprehensive experiments on various vision-language models to show that the proposed methods consistently improve over the baselines for different model families (*e.g.*, CLIP and AL-BEF) and for different downstream applications (*e.g.*, cross-modality retrieval, VQA, VR and *etc*).

## 2. Related work

Most recent works on vision-language representation learning can be categorized based on how information from different modalities is used for joint learning. The first category applies unified models [58, 59, 63, 64] to process both images and texts, where the inputs are usually tokenized into sequences [2, 47]. Unified models feature simpler and more universal designs, but typically underperform methods with modality-specific encoders (the second category). These methods use separate encoders for images and texts (*e.g.* CLIP [40, 48, 52], ALIGN [25]), and rely on contrastive loss [6, 20, 44] to align multiple modalities. These methods have been shown to achieve state-of-the-art (SOTA) performance on image-text retrieval; but the support is lacking for multi-modality tasks requiring inter-modality interaction, *e.g.* VQA. To conquer this problem, most recent approaches use a hybrid fashion where the models have separate encoders for images and texts along with a late-fusion multi-modal encoder [13, 30, 31, 53, 67]. Specifically, image-text matching (ITM) loss and masked language modeling (MLM) loss are usually applied for training the fusion encoder.

The methods in the later category utilize separate encoders for different modalities. However, this can lead to the phenomenon that image embeddings and text embeddings reside in different regions of the joint latent space. Such a phenomenon, termed *modality gap*, is observed in many multi-modal models [48, 66, 69]. A recent study [34] shows that the modality gap presents from the initialization and can be preserved during contrastive training. This naturally brings in another variety in multi-modality models – the latent modality gap and modality structures. Cy-CLIP [18] advocates for the benefit of consistency in latent modality structures. Yet to the best of our knowledge, no other prior work has studied the modality gap from a theoretical view. In this work, we show that directly reducing the modality gap does not help in performance gain from both empirical experiments and theoretical analysis. Consequently, we propose to study the impact of latent modality structures, and propose three approaches to obtain more meaningful latent modality structures that can improve downstream applications.

## 3. Understanding the Impact of Modality Gap on Downstream Performance

Despite being used extensively as a heuristic in practice [34, 66, 67, 69], it remains an open question whether modality alignment in the feature space through contrastive learning is optimal for downstream performance [34]. In this section, we first formally formulate the modality gap problem, present our empirical evidence on the relationship between the modality gap and the performance of down-

stream tasks, and then probe into its theoretical underpinning by providing an information-theoretical analysis.

**Notation** Throughout the paper, we will use $X_T$ and $X_V$ to denote the random variables corresponding to the input texts and images, respectively. We shall use $Y$ to denote the target variable in the downstream task of interest. For example, in the context of online shopping, $X_T$ and $X_V$ could be the textual and visual descriptions of a product, and in this case $Y$ is the expected sale of this product. When dealing with data with multi-modalities, we often use modality-specific encoder $g_T$ and $g_V$ to obtain features in the same latent space, i.e., $Z_T = g_T(X_T)$ and $Z_V = g_V(X_V)$ are the extracted features from textual and visual inputs. In this work, we focus on the setting where inputs from different modalities are paired with each other, meaning that a sample consists of the tuple $(x_T, x_V, y)$ from the underlying joint distribution $p$. The goal of reducing the modality gap in the latent space is then to shrink the statistical distance (e.g., KL-divergence, etc) between $Z_T$ and $Z_V$.

For two random variables $X_T$ and $X_V$, we define $I(X_T; X_V)$ to be the Shannon mutual information between $X_T$ and $X_V$. Similarly, we use $H(Y \mid X_T, X_V)$ to denote the conditional entropy of $Y$ given the two modalities as input. Following common practice, for classification tasks, $\ell_{\mathrm{CE}}(\hat{y}, y)$ is the cross-entropy loss between the prediction $\hat{y}$ and the ground-truth label $y$. One useful fact about the conditional entropy $H(Y \mid X_T, X_V)$ and the cross-entropy loss is the following variational form [14, 70]: $H(Y \mid X_T, X_V) = \inf_f \mathbb{E}_p[\ell_{\mathrm{CE}}(f(X_T, X_V), Y)]$, where the infimum is over all the prediction functions that take both $X_T$ and $X_V$ as input to predict the target $Y$ and the expectation is taken over the joint distribution $p$ of $(X_T, X_V, Y)$.

### 3.1. Empirical Analysis on Modality Gap

Given paired multi-modal data, one natural idea explored in the literature [34, 67, 69] is to use contrastive pretraining by treating paired multimodal data as the positive pairs and others as negative pairs. The goal is to align the positive pairs so that they are closer to each other in the feature space while at the same time ensuring the negative pairs to be farther away. More specifically, let $(x_T, x_V, y)$ and $(x_T', x_V', y')$ be two tuples sampled from the joint distribution. Then, in order to align the two modalities, $(x_T, x_V)$, $(x_T', x_V')$ are used as positive pairs while $(x_T, x_V')$ and $(x_T', x_V)$ are constructed as negative pairs.

Based on the contrastive loss principle [61, Theorem 1], a better model should come with smaller modality gaps (better alignment). However, despite being extensively used as a pretraining strategy in practice, it is unclear how the modality alignment affects the downstream tasks of interest. To approach this important question, we first conduct experiments to explore the effect of reducing modality gap
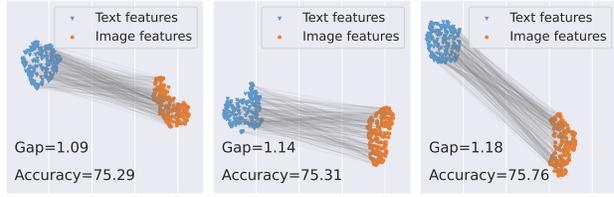


Figure 2. Visualization of the modality gap between text and image features. There is no clear-cut relationship between the gap of these two modalities and the downstream retrieval performance.

on the task of image/text retrieval.

We plot the alignment between paired image/text data in the feature space and also compute the average distance between them as the *gap* measure in Fig. 2. We perform pre-training on COCO [35] dataset and evaluate the zero-shot retrieval performance on Flick30K [68] test set. We optimize an additional alignment loss during training, $\mathcal{L}_{\mathrm{Align}} = 1/\langle Z_T, Z_V \rangle^2$, to reduce the gap between modalities. We control the gap by adjusting the scale of the alignment loss. From Fig. 2, we can see that the retrieval performance barely changes when changing the gap between two modalities. Note that as we normalized the data in the feature space, the gap difference in the figure is significant.

### 3.2. An Information-Theoretic Analysis on Modality Gap

Inspired by the empirical observation, we conjecture that *reducing the modality gap in feature space does not always lead to better downstream task performance*. Nevertheless, it is instructive to theoretically understand when and in what kind of downstream tasks reducing the modality gap could help. To do so, we first define the *information gap* $\Delta_p := |I(X_T; Y) - I(X_V; Y)|$ to characterize the gap of utility provided by two modalities towards predicting the target variable $Y$. Note that by definition, the information gap $\Delta_p$ only depends on the joint distribution $p$, i.e., the multimodal prediction problem itself, and is independent of the modality encoders $g_T$ and $g_V$. Hence, it is a constant during the modality learning process. As we shall see shortly, the *information gap* will serve as a lower bound of the downstream prediction error if we seek to find features that admit a *zero modality gap*. From this perspective, the information gap is the *price* we have to pay for using perfectly aligned features among different modalities. Thus, it well corresponds to the modality gap we are interested in. We can now state our theorem as follows.

**Theorem 3.1.** For a pair of modality encoders $g_T(\cdot)$ and $g_V(\cdot)$, if the multi-modal features $Z_T = g_T(X_T)$ and $Z_V = g_V(X_V)$ are perfectly aligned in the feature space, i.e., $Z_T = Z_V$, then $\inf_h \mathbb{E}_p[\ell_{\mathrm{CE}}(h(Z_T, Z_V), Y)] - \inf_{h'} \mathbb{E}_p[\ell_{\mathrm{CE}}(h'(X_T, X_V), Y)] \geq \Delta_p$.

**Remark** We discuss some of the implications of the above

theorem. At a high level, Theorem 3.1 states that if the information gap $\Delta_p$ between the two modalities is large, then the optimal prediction error we can hope to achieve by using modality-aligned features is at least $\Delta_p$ larger than that we can achieve from the input modalities. In particular, when only one of the modalities contains predictive information w.r.t. the downstream target $Y$, enforcing perfect modality alignment could render the learned modality-aligned features $Z_T$ and $Z_V$ uninformative of $Y$, leading to a large downstream prediction error. Intuitively, such a phenomenon will happen because modality alignment enforces the aligned features to only contain predictive information present in both of the input modalities $X_T$ and $X_V$.

In practice, because of the use of contrastive loss, due to the asymptotic behavior of it [61, Theorem 1], in the limit of infinity amount of data, the contrastive loss will force positive pairs to be perfectly aligned. In the context of multimodal learning, this means that the assumption $Z_T = Z_V$ of Theorem 3.1 will hold. As a last note, we comment that the requirement of perfect alignment in Theorem 3.1 is not necessary: the lower bound could be extended when the features $Z_T$ and $Z_V$ are only approximately aligned.[1]

Due to space limit, we defer the proof of Theorem 3.1 to Appendix A. In fact, it can be readily seen from the proof in the appendix that we could relax the exact modality alignment condition in Theorem 3.1 even further. In other words, as long as there exists a bijection between $Z_T$ and $Z_V$, then the conditional mutual information satisfies $I(Z_V; Y \mid Z_T) = I(Z_T; Y \mid Z_V) = 0$, so the exact same lower bound in Theorem 3.1 will hold.

## 4. Method

Motivated by Theorem 3.1, instead of seeking exact modality matching, in this section we propose to construct meaningful *latent modality structures*. They can play an important role in learning generalizable multi-modal representations by preventing pure modality alignment. In the following, we propose three designs from different perspectives to construct the latent modality structures, by considering variations in intra- and inter-modalities. We visualize these designs in Fig. 3. We first introduce the basic contrastive learning framework that we develop our methods on. Following previous work [13, 48], we adopt the multi-modal training framework with contrastive loss, which uses both cross-modal and in-modal contrastive loss,

*i.e.*, $\mathcal{L}_{\text{Con}} = \frac{1}{4}(\mathcal{L}_{\text{V2T}} + \mathcal{L}_{\text{T2V}} + \mathcal{L}_{\text{V2V}} + \mathcal{L}_{\text{T2T}})$ with:

$$\mathcal{L}_{\text{V2T}} = -\frac{1}{N} \sum_{j=1}^{N} \log \frac{e^{\langle z_{V_j}, z_{T_j} \rangle / \tau}}{\sum_{k=1}^{N} e^{\langle z_{V_j}, z_{T_k} \rangle / \tau}}$$

$$\mathcal{L}_{\text{V2V}} = -\frac{1}{N} \sum_{j=1}^{N} \log \frac{e^{\langle z_{V_j}, z_{V_j}^{\text{a}} \rangle / \tau}}{\sum_{k=1}^{N} e^{\langle z_{V_j}, z_{V_k} \rangle / \tau}}$$

where $N$ denotes the batch size; $z_{V_j}$ denote the feature of the $j$-th image in the mini-batch, with its augmentation $z_{V_j}^{\text{a}}$ and corresponding text feature $z_{T_j}$. The remaining losses ($\mathcal{L}_{\text{T2V}}$, $\mathcal{L}_{\text{T2T}}$) are defined in the same way by switching between text modality ($T$) and image modality ($V$).

### 4.1. Intra-modality Regularization via Deep Feature Separation
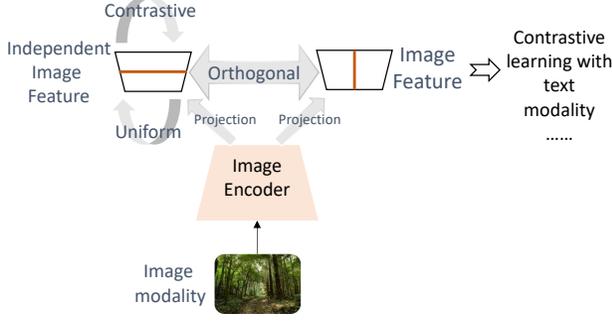
This subsection aims to construct intra-modality structures to regularize in-modality representations. Based on Theorem 3.1, we first define two types of information, *modality-shared information* that is shared by all modalities, and *modality-independent information* that is modality-specific. Our motivation stems from our theoretical finding that exact modality matching is sub-optimal due to the loss of *modality-independent information*. To overcome this limitation, we propose to explicitly model the modality-independent information. We achieve this by applying the idea of feature separation [4] on multi-modal representation learning. Our basic construction is shown in Figure 3a. On top of the contrastive learning framework, we use additional projection layers to construct new features to store such information. We term these *independent features*, meaning that they contain modality-specific information independent of the other modality. We take extra constraints to ensure that a) independent features contain complementary information from the original features; and b) independent features are meaningful representations.

To ensure a), we constrain the features to be orthogonal to the original features by forcing their inner product to be small, *i.e.* $\langle u, v \rangle = 0$. We define an orthogonal loss over minibatch optimization as follows:
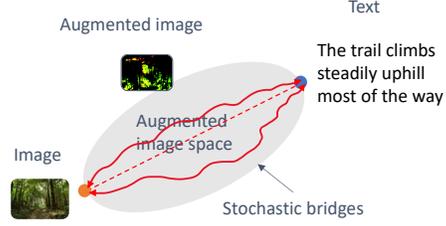
$$\mathcal{L}_{\text{Ortho}} = \frac{1}{N} \sum_{j=1}^{N} (\langle z_{V_j}, z_{V_j}^{\text{i}} \rangle^2 + \langle z_{T_j}, z_{T_j}^{\text{i}} \rangle^2)$$

where $z_{V_i}^{\text{i}}$ denote the independent feature of the $i^{th}$ image feature in the batch.
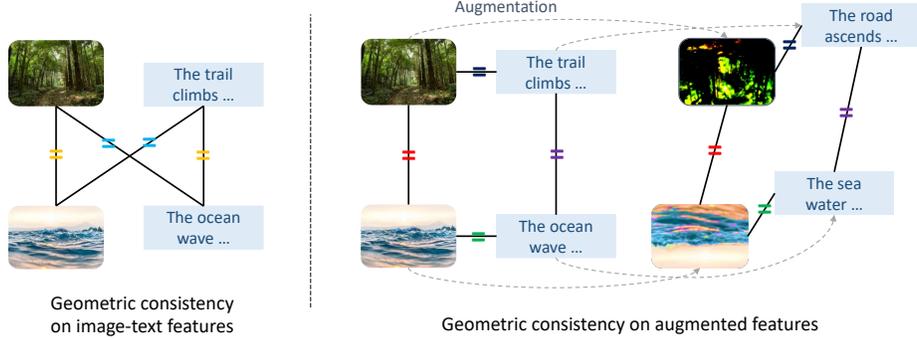
To avoid the degenerate case where the independent features are learned to be non-informative noises independent of the other modality, we further constrain that the independent features are informative. To this end, we adopt the contrastive loss and uniformity loss on the independent

(a) Building deep feature separation to preserve *modality-independent information*. Independent image features are enforced to be complimentary from original feature (with orthogonal loss) and store meaningful information (with contrastive and uniform losses).

(b) Building Brownian bridges between the image and text modalities to regularize inter-modality representations. Each red curve illustrates a stochastic bridge connecting an image-text pair; and the augmented images are enforced to stay on the path, guiding a cross-modality structure to connect the image and text modalities.

(c) Building geometric consistency between features. Each solid line represents the distance between two features. Same colored = signs indicate that the symmetry is encouraged, *i.e.* the two distances are encouraged to be the same.

Figure 3. Illustration of our three designed regularizer for constructing latent feature structure.

features, *i.e.*, we first adopt in-modality contrastive loss for independent text features and independent image features separately, *i.e.*, $\mathcal{L}_{\text{Con}}^{\text{i}} = \mathcal{L}_{\text{V2V}}^{\text{i}} + \mathcal{L}_{\text{T2T}}^{\text{i}}$ with

$$\mathcal{L}_{\text{V2V}}^{\text{i}} = -\frac{1}{N} \sum_{j=1}^{N} \log \frac{e^{\langle z_{V_j}^{\text{i}}, z_{V_j}^{\text{i a}} \rangle / \tau}}{\sum_{k=1}^{N} e^{\langle z_{V_j}^{\text{i}}, z_{V_k}^{\text{i}} \rangle / \tau}},$$

and $\mathcal{L}_{\text{T2T}}^{\text{i}}$ is defined similarly. Then we enhance the independent features with the uniformity loss [62] that maximizes the pairwise Gaussian potential [1,11]. Such a uniformity loss encourages the learned features to preserve maximal information:

$$\mathcal{L}_{\text{Uni}}^{\text{i}} = \log \frac{1}{N} \sum_{j=1}^{N} \sum_{k=1}^{N} G_t(z_{V_j}^{\text{i}}, z_{V_k}^{\text{i}}) + G_t(z_{T_j}^{\text{i}}, z_{T_k}^{\text{i}}),$$

where $G_t(u, v) = e^{-t\|u-v\|^2}$ is the Gaussian potential kernel with $t = 2$. In this way, we can preserve both *modality-shared information* and *modality-independent information*. Finally we obtain the total loss: $\mathcal{L}_{\text{Sep}} = \mathcal{L}_{\text{Ortho}} + \mathcal{L}_{\text{Con}}^{\text{i}} + \mathcal{L}_{\text{Uni}}^{\text{i}}$.

## 4.2. Inter-modality Regularization via Brownian Bridge

Next, we consider regularizing inter-modality structures. With the existence of modality gap, a natural idea is to con-

strain paired modality features in some subspace so that they are better separated from other feature pairs. To this end, we propose to construct a latent structure to explicitly guide the transition from the image modality to the associated text modality. Such a modality transition can be seamlessly modeled by the so-called Brownian bridge [39, 60], which is a special type of Brownian motion with constraints that define stochastic paths (called bridges) between a pair of fixed starting and ending points (corresponding to the two modalities in our setting). Our basic construction is illustrated in Figure 3b.

To formulate this, given two random variables $(Z_V, Z_T)$ of image-text feature pairs, we denote the feature of augmented image as $Z_V^{\text{a}}$. We define a stochastic path such that $Z_V^{\text{a}}$ is constrained to stay on the path between $Z_V$ and $Z_T$. From the property of Brownian bridge, this endows a conditional Gaussian distribution of the form:

$$p(Z_V^{\text{a}} | Z_V, Z_T) = \mathcal{N}(Z_V^{\text{a}}; \mu(Z_V, Z_T, t), t(1-t)\mathbf{I}) \quad (1)$$

where $t \in [0, 1]$ is a hyperparameter, which can be randomly sampled at each time or fixed to a pre-defined value (we fix it to 0.25 in our experiments for simplicity); $\mu(Z_V, Z_T, t) \triangleq \frac{t Z_V + (1-t) Z_T}{\|t Z_V + (1-t) Z_T\|}$, and the normalizer is applied in order to constrain the mean to lie on the hyper-

(a) Two-tower-based models.
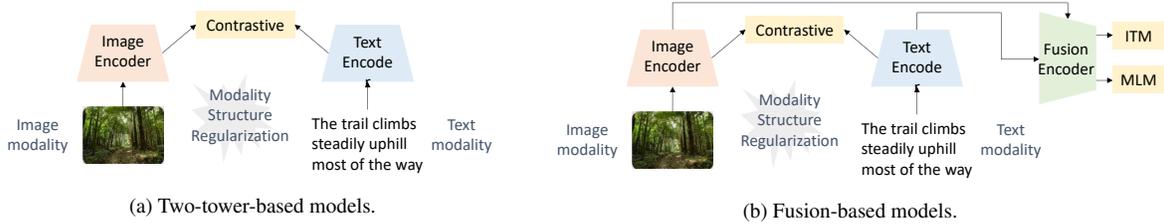
(b) Fusion-based models.

Figure 4. Illustration of two-tower-based models (*e.g.* CLIP) and fusion-based models (*e.g.* ALBEF). Our latent modality regularization can be applied to both type of models at their feature level.

sphere feature space. Based on the maximal likelihood principle, to fit the model, we can simply align the $Z_V^{\mathrm{a}}$ with the mean of the Brownian bridge in (1). When applying stochastic optimization, this ends up with optimizing the following objective at each time over a mini-batch:

$$\mathcal{L}_{\mathrm{Br}} = \frac{1}{N} \sum_{j=1}^{N} \| z_{V_j}^{\mathrm{a}} - \mu(z_{V_j}, z_{T_j}, t) \|^2$$

$$= \frac{1}{N} \sum_{j=1}^{N} \frac{t\langle z_{V_j}, z_{V_j}^{\mathrm{a}} \rangle + (1-t)\langle z_{T_j}, z_{V_j}^{\mathrm{a}} \rangle}{t^2 + (1-t)^2 + 2t(1-t)\langle z_{V_j}, z_{T_j} \rangle}$$

### 4.3. Intra-Inter Regularization via Geometric Consistency

In the previous subsections, we consider either intra- or inter-modality structures between the two modalities. Is it possible to relate these two types of relationships together? In this subsection, we aim to design a general regularizer that considers both intra- and inter-modality structures. We achieve this goal by enforcing geometric symmetry within and between modality representations and their augmentations. Specifically, we generalize the idea in CyCLIP [18] so that it also includes geometric consistency for the augmented features, which is demonstrated in the experiments to achieve significant improvement.

Specifically, we apply two types of geometric consistency losses that achieve symmetry in the following settings. First, we enforce geometric consistency among the original modality features, by optimizing the similarity between the mismatched image and text pairs, and the similarity between image pairs and text pairs. As shown in Figure 3c, we achieve this by encouraging the geometric consistency such that $\langle z_{V_1}, z_{T_2} \rangle \sim \langle z_{V_2}, z_{T_1} \rangle$ and $\langle z_{V_1}, z_{V_2} \rangle \sim \langle z_{T_1}, z_{T_2} \rangle$, where $a \sim b$ means $a$ is close to $b$ in some sense (defined below). We define the following geometric consistency objective over mini-batch:

$$\mathcal{L}_{\mathrm{GC}} = \frac{1}{N} \sum_{j=1}^{N} \sum_{k=1}^{N} [(\langle z_{V_j}, z_{T_k} \rangle - \langle z_{V_k}, z_{T_j} \rangle)^2$$

$$+ (\langle z_{V_j}, z_{V_k} \rangle - \langle z_{T_j}, z_{T_k} \rangle)^2]$$

Second, we optimize the geometric consistency of augmented features. As shown in Fig. 3c we optimize geometric symmetry between feature pairs and augmented feature

pairs in the text and image space. The following objective is used to enforce this goal:

$$\mathcal{L}_{\mathrm{GC}}^{\mathrm{a}} = \frac{1}{N} \sum_{j=1}^{N} \sum_{k=1}^{N} [(\langle z_{V_j}, z_{V_k} \rangle - \langle z_{V_j}^{\mathrm{a}}, z_{V_k}^{\mathrm{a}} \rangle)^2$$

$$+ (\langle z_{T_j}, z_{T_k} \rangle - \langle z_{T_j}^{\mathrm{a}}, z_{T_k}^{\mathrm{a}} \rangle)^2] + \frac{1}{N} \sum_{j=1}^{N} (\langle z_{V_j}, z_{T_j} \rangle - \langle z_{V_j}^{\mathrm{a}}, z_{T_j}^{\mathrm{a}} \rangle)^2$$

Overall, the total combination of geometric consistency loss can be written as: $\mathcal{L}_{\mathrm{GC}} + \mathcal{L}_{\mathrm{GC}}^{\mathrm{a}}$.

**Final Loss** We can now define a final loss by combining the standard contrastive loss with one or several of our proposed modality regularization losses. The effect of each regularization could be task-dependent, *i.e.* certain task could benefit more from certain regularization, which we will show comprehensively in the next section.

## 5. Experiments

Our proposed methods are general purposed. Thus, we choose to evaluate them with two popular multi-modal representation frameworks: the two-tower based models (*e.g*, CLIP) and the fusion based models (*e.g.*, ALBEF), as illustrated in Fig. 4. Note that in CLIP, text inputs are augmented with EDA [65], and image inputs are augmented with random augmentation such as flipping and cropping. In ALBEF, augmented features are obtained with additional momentum encoders.

### 5.1. Two-Tower-based Models

For this set of experiments, we adopt the CLIP-based models, where two separate encoders are trained to align features from the image and text modalities. To regularize latent modality structures, our regularization losses are separately applied along with the standard contrastive loss for pre-training[2]. We then evaluate on standard benchmarks.

**Setup:** Our CLIP model adopts ResNet-50 [21] as the image encoder and BERT [12] as the text encoder. We adopt

---

[2]We will combine all the proposed regularizers for evaluation in experiments with the fusion-based models.

Table 1. Zero-shot TopK classification accuracy (%) on CIFAR10, CIFAR100 and ImageNet1K.

| Method | CIFAR10 | | | CIFAR100 | | | ImageNet1K | | |
|---|---|---|---|---|---|---|---|---|---|
| | Top1 | Top3 | Top5 | Top1 | Top3 | Top5 | Top1 | Top3 | Top5 |
| CLIP [48] | 44.95 | 72.58 | 88.3 | 15.05 | 29.51 | 37.53 | 16.72 | 28.61 | 34.38 |
| CyCLIP [18] | 43.22 | 71.43 | 83.22 | 15.09 | 27.39 | 34.35 | 17.77 | 30.06 | 36.20 |
| OURS$_{Sep}$ | 46.61 | **81.21** | **92.44** | 19.37 | 36.66 | 46.26 | 20.21 | 33.25 | 39.60 |
| OURS$_{Br}$ | 43.15 | 72.77 | 86.72 | 14.22 | 26.46 | 33.28 | **20.45** | **33.56** | 39.28 |
| OURS$_{GC}$ | **56.36** | 80.47 | 90.27 | **22.70** | **41.66** | **51.78** | 20.25 | 33.50 | **39.91** |

Table 2. Zero-shot TopK classification accuracy (%) on Natural Distribution Shifts.

| Method | ImageNetV2 | | | ImageNetSketch | | | ImageNet-A | | | ImageNet-R | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Top1 | Top3 | Top5 | Top1 | Top3 | Top5 | Top1 | Top3 | Top5 | Top1 | Top3 | Top5 |
| CLIP [48] | 14.11 | 25.76 | 31.80 | 8.61 | 16.47 | 21.13 | 2.81 | 7.31 | 11.32 | 19.07 | 31.99 | 39.03 |
| CyCLIP [18] | 15.25 | 26.59 | 32.15 | 8.30 | 16.18 | 20.77 | 3.27 | 8.45 | 13.07 | 19.85 | 33.35 | 40.35 |
| OURS$_{Sep}$ | 16.78 | 28.97 | 35.68 | 9.22 | 17.86 | 23.00 | 3.45 | 9.88 | 15.81 | 22.06 | 35.65 | 43.01 |
| OURS$_{Br}$ | 17.02 | 29.39 | 35.53 | 10.34 | 18.39 | 23.05 | 3.01 | 7.50 | 11.45 | 20.40 | 32.43 | 38.45 |
| OURS$_{GC}$ | **17.37** | **29.84** | **36.65** | **10.90** | **20.77** | **26.11** | **3.87** | **11.36** | **16.76** | **23.85** | **37.90** | **45.03** |

Table 3. Linear probing Top1 classification accuracy (%) on visual benchmarks.

| | Caltech101 | SVHN | STL10 | CIFAR10 | CIFAR100 | DTD | FGVCAircraft | OxfordPets | SST2 | Food101 | GTSRB | StanfordCars | Flowers102 | ImageNet1K | Average |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CLIP [48] | 78.57 | 57.07 | 87.22 | 79.74 | 56.36 | 59.84 | 37.17 | 59.66 | 53.98 | 58.11 | 74.21 | 23.96 | 76.66 | 52.10 | 61.05 |
| CyCLIP [18] | 77.86 | 54.29 | 87.61 | 77.53 | 54.23 | 58.19 | 33.00 | 62.63 | 54.81 | 60.82 | 72.95 | 23.36 | 72.89 | 52.83 | 60.14 |
| OURS$_{Sep}$ | **84.45** | **69.82** | 90.96 | **81.51** | **61.19** | **67.50** | **41.70** | 67.16 | 54.26 | **63.08** | **82.35** | **31.76** | **81.69** | **56.73** | **66.73** |
| OURS$_{Br}$ | 82.18 | 57.46 | 90.69 | 79.42 | 57.72 | 64.84 | 34.74 | 65.71 | 54.04 | 60.52 | 73.61 | 26.50 | 78.44 | 53.87 | 62.84 |
| OURS$_{GC}$ | 83.23 | 63.58 | **91.31** | 80.92 | 58.89 | 65.43 | 34.83 | 64.51 | **55.19** | 60.80 | 76.84 | 26.95 | 78.76 | 54.96 | 64.01 |

the official code from CyCLIP to incorporate our regularizations, as well as to reproduce the baselines. Our reproduced CLIP results are consistent with the recent works [17, 41], although they are slightly lower than reported in the original CLIP paper. The reason could be that the number of GPUs we use is different and we provide details in Appendix C.1. For both baselines, we can reproduce better performance on linear probing but slightly under-perform on zero-shot transfer, which we consider reasonable. Note that all methods are under the same codebase and same hyper-parameter setting, thus the comparisons are fair.

**Pre-training:** We follow the protocol of previous works to pre-train the model with the CC3M [51] dataset, which contains 3M unique images and 4M image-text pairs.

### 5.1.1 Zero-Shot Transfer Learning Evaluation

We perform zero-shot transfer on standard image classification tasks, with the CIFAR10, CIFAR100 [29] and ImageNet1K [50] datasets. We use the standard evaluation strategy of prompt engineering. For each dataset, we construct the text prompts using the name of the class, *e.g.* "a photo of the [class name]". For each class, we obtain the normalized class text embedding. During the evaluation, the class with the highest similarity score to the im-

age embedding is predicted to be the label. Following previous works, we report Top-K classification accuracy with $K = 1, 3, 5$.

As shown in Tab. 1, our method significantly outperforms CLIP and CyCLIP on all three datasets, demonstrating the importance of latent modality structures. It is also interesting to see the differences our three regularizers perform in different datasets, *i.e.*, the feature-separation regularizer performs best in CIFAR10, while Brownian bridge regularizer performs best on ImageNet1K, and geometry consistency regularizer performs the best on CIFAR100.

### 5.1.2 Natural Distribution Shift Evaluation

We further evaluate variants [22, 23, 49, 57] of ImageNet1K dataset with shifted distributions. These datasets contain sketches, cartoons and adversarial generated images. As shown in Tab. 2, all methods suffer from performance degradation on natural distribution shift benchmarks compared to the performance on original ImageNet1K in Tab. 1. Nevertheless, our method consistently outperforms the baselines on all benchmarks. In contrast to the other experiments, our geometric consistency regularization performs the best on all the benchmarks.

Table 4. Downstream tasks performance on fusion-based models.

| Method | VQA | | NLVR$^2$ | | SNLI-VE | |
|---|---|---|---|---|---|---|
| | test-dev | test-std | dev | test-P | val | test |
| ImageBERT [32] | 70.80 | 71.00 | 67.40 | 67.00 | - | - |
| LXMERT [56] | 72.42 | 72.54 | 74.90 | 74.50 | - | - |
| 12-in-1 [37] | 73.15 | - | - 78.87 | - | 76.95 | |
| UNITER [7] | 72.70 | 72.91 | 77.81 | 77.85 | 78.59 | 78.28 |
| OSCAR [33] | 73.16 | 73.44 | 78.07 | 78.36 | - | - |
| VILLA [16] | 73.59 | 73.67 | 78.39 | 79.30 | 79.47 | 79.03 |
| ViLT [26] | 70.94 | - | 75.24 | 76.21 | - | - |
| ViCHA [53] | 73.55 | - | 78.14 | 77.00 | 79.20 | 78.65 |
| ALBEF [31] | 73.38 | 73.52 | 78.36 | 79.54 | 79.69 | 79.91 |
| CODIS [13] | 73.15 | 73.29 | 78.58 | **79.92** | 79.45 | 80.13 |
| OURS$_{Sep}$ | 73.52 | 73.59 | **79.05** | 79.76 | **79.95** | 79.61 |
| OURS$_{Br}$ | **74.26** | **74.36** | 78.70 | 79.36 | 79.86 | 79.95 |
| OURS$_{GC}$ | 73.90 | 73.87 | 78.96 | 79.53 | 79.82 | **80.16** |

### 5.1.3 Linear Probing Evaluation

We demonstrate better latent structure can also benefit downstream tasks with in-domain supervision. We evaluate this on linear probing tasks by fitting a linear classifier with in-domain supervision using the learned visual encoder. In total, we evaluate on 14 standard benchmarks [3,9,10,15,24,27,29,38,42,43,46,50,54]. As shown in Tab. 3, all our methods outperform the baselines on all benchmarks by large margins. Remarkably, our deep feature separation regularization performs particularly well on this task. We believe this is partially because such regularization can learn to preserve more information that could be useful with extra in-domain supervision.

## 5.2. Fusion-based Models

We next test our methods on fusion-based models. We adopt the ALBEF [31] framework, where a fusion encoder is applied to fuse the modality as shown in Fig. 7b. Such fusion-based models are known to be more powerful in learning inter-model interaction compared to simple two-tower-based models. Thus, we evaluate our methods on various vision-language downstream tasks including VQA [19], NLVR$^2$ [55], SNLI-VE [5]. Here we incorporate all three regularizations for these tasks. We additionally provide ablation study on smaller scale experiments.

**Setup** We use ViT-B/16 as our vision encoder and 12-layer BERT$_{base}$ as the text encoder. Note the first 6 layers of BERT$_{base}$ are used purely as the text encoder and the remaining are used as fusion encoder. We reproduced ALBEF and CODIS results for fair comparisons. All experiments we run are under the same codebase and hyper-parameter settings. The details are included in Appendix C.2.

**Pre-training:** We follow the previous experiments protocols [13,31] using a union of four datasets for pre-training, which include Conceptual Captions (CC3M) [51], Visual Genome (VG) [28], SBU Captions [45] and COCO [35], constituting 4M unique images and 5M image-text pairs.

### 5.2.1 Vision-Language Tasks Evaluation

**Visual Question Answering (VQA):** We fine-tune and evaluate our pre-trained model on VQA v2.0. Following [8,13,31], we consider VQA as a generation task. During fine-tuning, we apply 6-layer transformer-based decoder to generate the answer. We fine-tune on the training set and evaluate on the test-dev and test-std set. The results are presented in Table 4. Consistently, our method performs the best and achieves a 1% improvement on both the test-dev and test-std sets.

**Natural Language for Visual Reasoning (NLVR$^2$):** We use the NLVR$^2$ dataset, which contains 100K texts paired with web images. To enable our model to reason over two images, we follow [31] to extend the fusion encoder with an MLP prediction head and perform additional pre-training of one epoch to prepare the fusion encoder on text-assignment task. As shown in Table 4, our method achieves an improvement of 2% on the dev set and matches the performance of SOTA on the test-P set.

**Visual Entailment (VE):** We follow [7,31] and consider this as a classification problem with three classes (entailment, neutral, contradictory). Thus, we adopt an MLP prediction head on top of the fusion encoder. Again, our method is comparable to the baselines on the val set and outperforms all baselines on the test set.

We provide additional results including analysis and visualization of constructing latent structures, visualization of experimental results, as well as ablation studies in Appendix B.

## 6. Conclusion

In this paper, we investigate the latent modality structures in multi-modal representation learning. We analyze and examine the modality gap in the latent feature space and reveal that reducing modality gap to zero does not always lead to better performance. Instead we advocate that more meaningful latent features structures will benefit the downstream applications. Thus we design three regularization methods to construct meaningful latent structures. We propose to use 1) deep feature separation loss 2) brownian bridge loss 3) geometric consistency loss to improve the latent features from different perspectives. Extensive experiments on multiple vision-language tasks including image classification, linear probing, visual question answering, visual reasoning, visual entailment confirm the effectiveness and the generalizability of our proposed approach on popular contrastive representation learning frameworks.

# References

[1] Philip Bachman, R. Devon Hjelm, and William Buchwalter. Learning representations by maximizing mutual information across views. In *Proc. NeurIPS*, 2019. 5

[2] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *ArXiv*, abs/2106.08254, 2022. 2

[3] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative components with random forests. In *Proc. ECCV*, 2014. 8

[4] Konstantinos Bousmalis, George Trigeorgis, Nathan Silberman, Dilip Krishnan, and Dumitru Erhan. Domain separation networks. In *Proc. NeurIPS*, 2016. 4

[5] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal, Sept. 2015. Association for Computational Linguistics. 8

[6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *Proc. ICML*, 2020. 2

[7] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *Proc. ECCV*, 2020. 8

[8] Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. Unifying vision-and-language tasks via text generation. In *Proc. ICML*. PMLR, 2021. 8

[9] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. Describing textures in the wild. In *Proc. CVPR*, 2014. 8

[10] Adam Coates, A. Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *AISTATS*, 2011. 8

[11] Henry Cohn and Abhinav Kumar. Universally optimal distribution of points on spheres. *arXiv: Metric Geometry*, 2006. 5

[12] 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. 6

[13] Jiali Duan, Liqun Chen, Son Tran, Jinyu Yang, Yi Xu, Belinda Zeng, and Trishul Chilimbi. Multi-modal alignment using representation codebook. In *Proc. CVPR*, 2022. 1, 2, 4, 8, 12, 13

[14] Farzan Farnia and David Tse. A minimax approach to supervised learning. In *Proc. NeurIPS*, volume 29, 2016. 3

[15] Li Fei-Fei, R. Fergus, and P. Perona. One-shot learning of object categories. *IEEE TPAMI*, 2006. 8

[16] Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu, Yu Cheng, and Jingjing Liu. Large-scale adversarial training for vision-and-language representation learning. *ArXiv*, abs/2006.06195, 2020. 8

[17] Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. Clip-adapter: Better vision-language models with feature adapters. *arXiv preprint arXiv:2110.04544*, 2021. 7

[18] Shashank Goel, Hritik Bansal, Sumit Kaur Bhatia, Ryan A. Rossi, Vishwa Vinay, and Aditya Grover. Cyclip: Cyclic contrastive language-image pretraining. *ArXiv*, abs/2205.14459, 2022. 2, 6, 7, 12

[19] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 8

[20] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proc. CVPR*, 2020. 2

[21] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *Proc. CVPR*, 2016. 6

[22] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Lixuan Zhu, Samyak Parajuli, Mike Guo, Dawn Xiaodong Song, Jacob Steinhardt, and Justin Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proc. ICCV*, 2021. 7

[23] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Xiaodong Song. Natural adversarial examples. In *Proc. CVPR*, 2021. 7

[24] Sebastian Houben, Johannes Stallkamp, Jan Salmen, Marc Schlipsing, and Christian Igel. Detection of traffic signs in real-world images: The German Traffic Sign Detection Benchmark. In *International Joint Conference on Neural Networks*, number 1288, 2013. 8

[25] 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 *Proc. ICML*, 2021. 1, 2

[26] Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision. In *Proc. ICML*, 2021. 8

[27] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *IEEE Workshop on 3D Representation and Recognition (3dRR-13)*, 2013. 8

[28] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *IJCV*, 123(1), 2017. 8

[29] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009. 7, 8

[30] Gukyeong Kwon, Zhaowei Cai, Avinash Ravichandran, Erhan Bas, Rahul Bhotika, and Stefan 0 Soatto. Masked vision and language modeling for multi-modal representation learning. *ArXiv*, abs/2208.02131, 2022. 1, 2

[31] Junnan Li, Ramprasaath R. Selvaraju, Akhilesh Deepak Gotmare, Shafiq Joty, Caiming Xiong, and Steven Hoi. Align before fuse: Vision and language representation learning with momentum distillation. In *Proc. NeurIPS*, 2021. 1, 2, 8, 12, 13

[32] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *ArXiv*, abs/1908.03557, 2019. 8

[33] Xiujun Li, Xi Yin, Chunyuan Li, Xiaowei Hu, Pengchuan Zhang, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng Gao. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *Proc. ECCV*, 2020. 8

[34] Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Zou. Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning. In *Proc. NeurIPS*, 2022. 1, 2, 3

[35] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Proc. ECCV*, 2014. 3, 8

[36] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *Proc. ICLR*, 2019. 13

[37] Jiasen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi Parikh, and Stefan Lee. 12-in-1: Multi-task vision and language representation learning. In *Proc. CVPR*, 2020. 8

[38] S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi. Fine-grained visual classification of aircraft. Technical report, 2013. 8

[39] Roger Mansuy and Marc Yor. Aspects of brownian motion. 2008. 5

[40] Ron Mokady, Amir Hertz, and Amit H Bermano. Clipcap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*, 2021. 1, 2

[41] Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. Slip: Self-supervision meets language-image pretraining, 2021. 7

[42] Yuval Netzer, Tao Wang, Adam Coates, A. Bissacco, Bo Wu, and A. Ng. Reading digits in natural images with unsupervised feature learning. In *NeurIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011. 8

[43] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. *Indian Conference on Computer Vision, Graphics & Image Processing*, 2008. 8

[44] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 2

[45] Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million captioned photographs. In *Proc. NeurIPS*, 2011. 8

[46] Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In *Proc. CVPR*, 2012. 8

[47] Zhiliang Peng, Li Dong, Hangbo Bao, Qixiang Ye, and Furu Wei. Beit v2: Masked image modeling with vector-quantized visual tokenizers. *ArXiv*, abs/2208.06366, 2022. 2

[48] 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 *Proc. ICML*, 2021. 1, 2, 4, 7

[49] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? In *Proc. ICML*, 2019. 7

[50] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115, 2015. 7, 8

[51] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, 2018. 7, 8

[52] Sheng Shen, Liunian Harold Li, Hao Tan, Mohit Bansal, Anna Rohrbach, Kai-Wei Chang, Zhewei Yao, and Kurt Keutzer. How much can clip benefit vision-and-language tasks? *arXiv preprint arXiv:2107.06383*, 2021. 1, 2

[53] Mustafa Shukor, Guillaume Couairon, and Matthieu Cord. Efficient vision-language pretraining with visual concepts and hierarchical alignment. *ArXiv*, abs/2208.13628, 2022. 1, 2, 8

[54] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proc. EMNLP*. Association for Computational Linguistics, 2013. 8

[55] Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2019. 8

[56] Hao Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. *ArXiv*, abs/1908.07490, 2019. 8

[57] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. In *Proc. NeurIPS*, 2019. 7

[58] Jianfeng Wang, Xiaowei Hu, Zhe Gan, Zhengyuan Yang, Xiyang Dai, Zicheng Liu, Yumao Lu, and Lijuan Wang. Ufo: A unified transformer for vision-language representation learning. *ArXiv*, abs/2111.10023, 2021. 1, 2

[59] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *Proc. ICML*, 2022. 1, 2

[60] Rose E. Wang, Esin Durmus, Noah D. Goodman, and Tatsunori Hashimoto. Language modeling via stochastic processes. *ArXiv*, abs/2203.11370, 2022. 5

[61] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*, pages 9929–9939. PMLR, 2020. 3, 4

[62] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *Proc. ICML*. PMLR, 2020. 5

[63] Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Mohammed, Saksham Singhal, Subhojit Som, and Furu Wei. Image as a foreign language: Beit pretraining for all vision and vision-language tasks. *ArXiv*, abs/2208.10442, 2022. 1, 2

[64] Wenhui Wang, Hangbo Bao, Li Dong, and Furu Wei. Vlmo: Unified vision-language pre-training with mixture-of-modality-experts. *ArXiv*, abs/2111.02358, 2021. 1, 2

[65] Jason Wei and Kai Zou. EDA: Easy data augmentation techniques for boosting performance on text classification tasks. In *Proc. EMNLP*. Association for Computational Linguistics, 2019. 6

[66] Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and Christoph Feichtenhofer. VideoCLIP: Contrastive pre-training for zero-shot video-text understanding. In *Proc. EMNLP*. Association for Computational Linguistics, 2021. 1, 2

[67] Jinyu Yang, Jiali Duan, Son Tran, Yi Xu, Sampath Chanda, Liqun Chen, Belinda Zeng, Trishul Chilimbi, and Junzhou Huang. Vision-language pre-training with triple contrastive learning. In *Proc. CVPR*, 2022. 1, 2, 3

[68] Peter Young, Alice Lai, Micah Hodosh, and J. Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2, 2014. 3, 12

[69] Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis P Langlotz. Contrastive learning of medical visual representations from paired images and text. *arXiv preprint arXiv:2010.00747*, 2020. 1, 2, 3

[70] Han Zhao, Chen Dan, Bryon Aragam, Tommi S. Jaakkola, Geoffrey J. Gordon, and Pradeep Ravikumar. Fundamental limits and tradeoffs in invariant representation learning, 2020. 3

# A. Proof of Theorem 3.1

To ease the reading, we first restate Theorem 3.1 and then provide the proof.

**Theorem 3.1.** For a pair of modality encoders $g_T(\cdot)$ and $g_V(\cdot)$, if the multi-modal features $Z_T = g_T(X_T)$ and $Z_V = g_V(X_V)$ are perfectly aligned in the feature space, i.e., $Z_T = Z_V$, then $\inf_h \mathbb{E}_p[\ell_{CE}(h(Z_T, Z_V), Y)] - \inf_{h'} \mathbb{E}_p[\ell_{CE}(h'(X_T, X_V), Y)] \geq \Delta_p$.

*Proof of Theorem 3.1.* Consider the joint mutual information $I(Z_T, Z_V; Y)$. By the chain rule, we have the following decompositions:

$$I(Z_T, Z_V; Y) = I(Z_T; Y) + I(Z_V; Y \mid Z_T)$$
$$= I(Z_V; Y) + I(Z_T; Y \mid Z_V).$$

However, since $Z_T$ and $Z_V$ are perfectly aligned, $I(Z_V; Y \mid Z_T) = I(Z_T; Y \mid Z_V) = 0$, which means $I(Z_T, Z_V; Y) = I(Z_V; Y) = I(Z_T; Y)$. On the other hand, by the celebrated data-processing inequality, we know that

$$I(Z_T; Y) \leq I(X_T; Y), \quad I(Z_V; Y) \leq I(X_V; Y).$$

Hence, the following chain of inequalities holds:

$$I(Z_T, Z_V; Y) = \min\{I(Z_T; Y), I(Z_V; Y)\}$$
$$\leq \min\{I(X_T; Y), I(X_V; Y)\}$$
$$\leq \max\{I(X_T; Y), I(X_V; Y)\}$$
$$\leq I(X_T, X_V; Y),$$

where the last inequality follows from the fact that the joint mutual information $I(X_T, X_V; Y)$ is at least as large as any one of $I(X_T; Y)$ and $I(X_V; Y)$. Therefore, due to the variational form of the conditional entropy, we have

$$\inf_h \mathbb{E}_p[\ell_{CE}(h(Z_T, Z_V), Y)] - \inf_{h'} \mathbb{E}_p[\ell_{CE}(h'(X_T, X_V), Y)]$$
$$= H(Y \mid Z_T, Z_V) - H(Y \mid X_T, X_V)$$
$$= I(X_T, X_V; Y) - I(Z_T, Z_V; Y)$$
$$\geq \max\{I(X_T; Y), I(X_V; Y)\} - \min\{I(X_T; Y), I(X_V; Y)\}$$
$$= \Delta_p. \qquad \blacksquare$$

# B. Additional Results

## B.1. Two-tower-based models

**Visualization of constructing latent structures** To better understand the effect of constructing latent modality structures, we visualize the effect of our method on the latent space in Fig. 5. Note that all our methods achieve performance gain regardless of size of the modality gap, which complys with Section 3 and Theorem 3.1.

Table 5. Downstream tasks performance on fusion-based models.

| Method | VQA | | NLVR$^2$ | | SNLI-VE | |
|---|---|---|---|---|---|---|
| | test-dev | test-std | dev | test-P | val | test |
| ALBEF [31] | 73.38 | 73.52 | 78.36 | 79.54 | 79.69 | 79.91 |
| CODIS [13] | 73.15 | 73.29 | 78.58 | **79.92** | 79.45 | 80.13 |
| OURS$_{All}$ | 74.12 | 74.16 | **80.18** | 79.80 | 79.62 | **80.23** |
| OURS$_{Sep}$ | 73.52 | 73.59 | 79.05 | 79.76 | **79.95** | 79.61 |
| OURS$_{Br}$ | **74.26** | **74.36** | 78.70 | 79.36 | 79.86 | 79.95 |
| OURS$_{GC}$ | 73.90 | 73.87 | 78.96 | 79.53 | 79.82 | 80.16 |

Table 6. Ablation study on zero-shot image-text retrieval performance on Flickr30K with model pre-trained on COCO.

| Method | Text Retrieval | | | Image Retrieval | | |
|---|---|---|---|---|---|---|
| | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 |
| ALBEF [31] | 58.4 | 83.2 | 89.5 | 44.5 | 69.8 | 78.0 |
| CODIS [13] | 62.7 | 87.0 | 92.3 | 49.0 | 74.1 | 82.9 |
| OURS$_{Sep}$ | **66.0** | **88.2** | **93.9** | 50.4 | 76.2 | 83.7 |
| OURS$_{Br}$ | 65.4 | 88.1 | 93.1 | **50.8** | **77.1** | **84.4** |
| OURS$_{GC}$ | 64.3 | 87.5 | 92.3 | 50.5 | 75.9 | 83.3 |

**Visualization of experimental results** To better demonstrate the effectiveness of our proposed methods, we visualize our experimental results on two-tower-based framework (*e.g.* CLIP) in Fig. 6 and Fig. 7. Our methods show significant improvement on most of the tasks.

## B.2. Fusion-based models

**Additional experimental results** We provide additional results on using all three regularizers. The results are shown in Tab. 5. While using all the regularizations together leads to performance gain, all our regularization methods improve the performance as well when used individually.

We evaluate zero-shot image-text retrieval on smaller scale experiments by pertaining on COCO and evaluate on Flickr30 [68]. As shown in Tab. 6, results indicate that all three regularizations improve the performance, while text retrieval benefits most from deep feature separation regularization and image retrieval task benefits most from Brownian bridge regularization.

# C. Implementation Details

## C.1. Two-tower based models.

We follow the same code base and hyper-parameters setting as CyCLIP [18] except for number of GPUs. We train the model from scratch on 64 NVIDIA A100 GPUs and train for 64 epochs. Our batch size is 128 and feature dimension is 1024. We use an initial learning rate of $5e^{-4}$ with cosine scheduling. We warm-up the model for 10000 steps. We evaluate the model trained to the last epoch for our method.

## C.2. Fusion based models.

We follow the codebase and hyper-parameter setting as [13, 31] except for number of GPUs. We train all the models on 16 NVIDIA A100 GPUs. During the pre-train stage, we train with the pre-training tasks for 30 epochs. AdamW [36] optimizer is used along with weight decay of 0.02, batch size of 512, learning rate initially of 1e-5. We warm up the learning rate to 1e-4 after 1000 iterations and follow the cosine decay. The input size for pre-training task is 256 and the input sizes for downstream tasks are 384.

**Reproducibility** We follow the standard practice to fix the random seed to ensure that all our results are reproducible. All the source code will be made public upon acceptance of the paper.
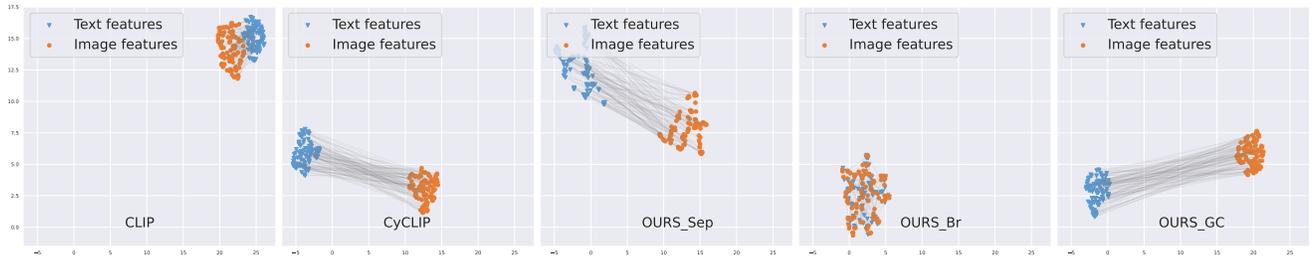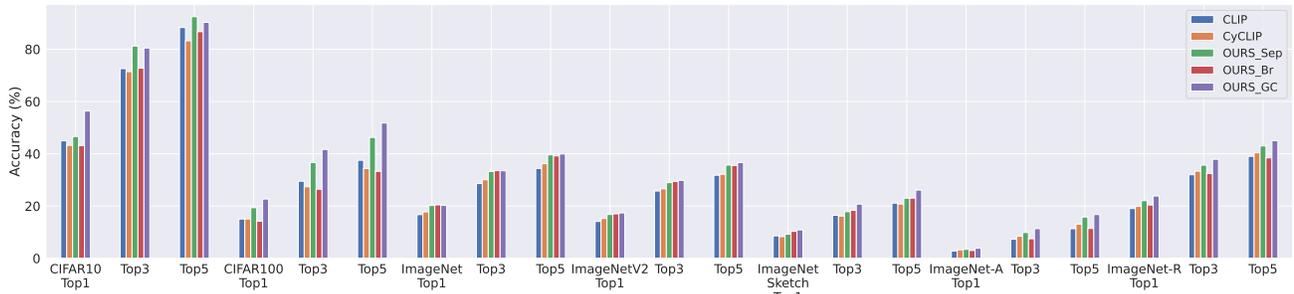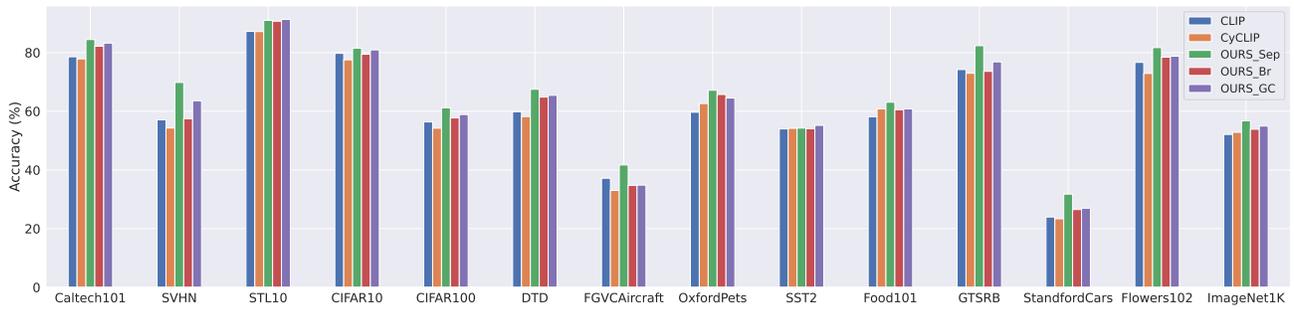
Figure 5. Visualization of constructing latent modality structures. Each line connects the positive image-text feature pair.
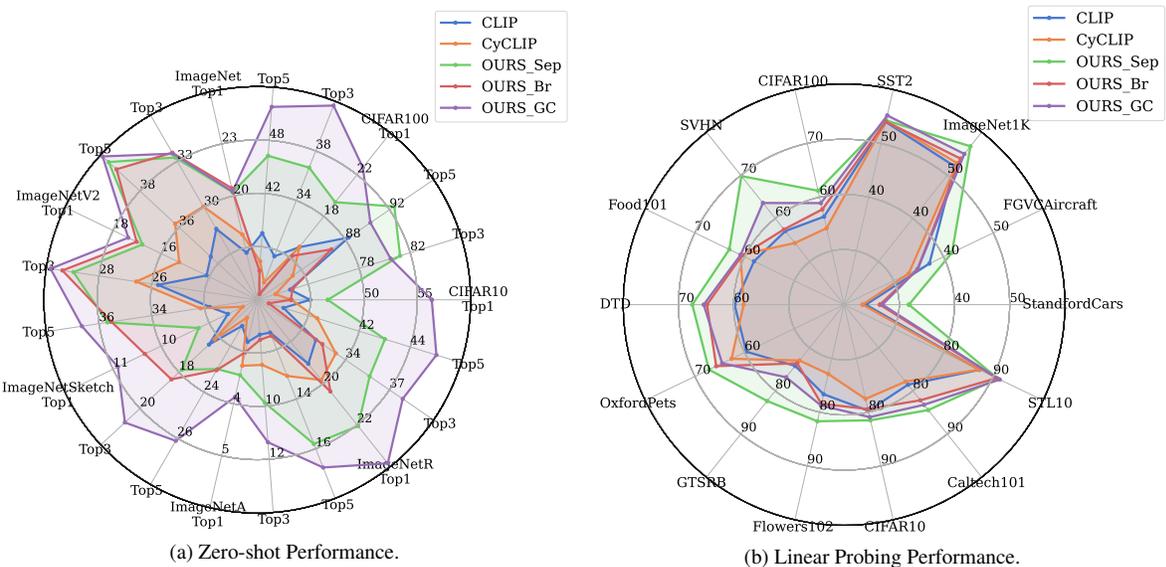


(a) Zero-shot transfer performance.



(b) Linear-probing performance.

Figure 6. Visualization of Two-tower-based methods (*e.g.* CLIP) performance. Each color represents a different approach.



(a) Zero-shot Performance.

(b) Linear Probing Performance.

Figure 7. Visualization of Two-tower-based methods (*e.g.* CLIP) performance. Each axis represents the performance on a dataset with a certain metric. Each color represents different approach. The larger area that one approach covers, the better overall performance.