

ANCHORS AWEIGH! SAIL FOR OPTIMAL UNIFIED MULTI-MODAL REPRESENTATIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Multimodal learning plays a crucial role in enabling machine learning models to fuse and utilize diverse data sources, such as text, images, and audio, to support a variety of downstream tasks. A unified representation across various modalities is particularly important for improving efficiency and performance. Recent binding methods, such as ImageBind (Girdhar et al., 2023), typically use a fixed anchor modality to align multimodal data in the anchor modal embedding space. In this paper, we mathematically analyze the *fixed anchor binding methods* and uncover notable limitations: (1) over-reliance on the choice of the anchor modality, (2) failure to capture intra-modal information, and (3) failure to account for inter-modal correlation among non-anchored modalities. To address these limitations, we propose CentroBind, a simple yet powerful approach that eliminates the need for a fixed anchor; instead, it employs dynamically adjustable centroid-based anchors generated from all available modalities, resulting in a balanced and rich representation space. We theoretically demonstrate that our method captures three crucial properties of multimodal learning: intra-modal learning, inter-modal learning, and multimodal alignment, while also constructing a robust unified representation across all modalities. Our experiments on both synthetic and real-world datasets demonstrate the superiority of the proposed method, showing that dynamic anchor methods outperform all fixed anchor binding methods as the former captures more nuanced multimodal interactions.

1 INTRODUCTION

Multimodal alignment is defined as identifying and exploiting relationships and correspondences between sub-components of instances from multiple modalities (e.g., text, image, audio) to establish meaningful connections between their representations (Baltrušaitis et al., 2018). This process allows machine learning models to analyze heterogeneous data holistically, facilitating comprehensive decision-making. A common approach is learning a shared embedding space (Tu et al., 2022; Girdhar et al., 2023; Liang et al., 2024b; Zhu et al., 2024), which aims to project data from multiple modalities into a common embedding space by clustering similar items together for direct comparison and linkage. This approach leverages well-trained single-modal embeddings, aligning them with auxiliary objective functions like contrastive (Oord et al., 2018) or triplet loss (Wang et al., 2020b) to minimize distances between similar items and maximize distances between dissimilar ones across modalities.

Instead of training separate models for each modality, ImageBind (Girdhar et al., 2023) pairs **images** with other modalities and projects them into the common image embedding space. Similarly, Zhu et al. (2024) shows that pairing **texts** with other modalities (LanguageBind) improves retrieval task performance when language is specified as the anchored modality. This approach has inspired various “-Bind” methods tailored to align different modalities for specific domains, such as molecular modeling (Xiao et al., 2024), medical imaging (Gao et al., 2024), brain signals (Yang et al., 2024b), and music selection for videos (Teng et al., 2024). These models commonly use image or text as the **anchor embedding** due to the abundance of data, with other modalities projected into this anchor representation.

We can define the aforementioned approaches as **Fixed-Anchor-Bind** (FABIND) method, where the embedding space of the primary anchor modality remains fixed during the alignment process. Generally, the “-Bind”-like approaches maximize mutual information $I(\mathbf{Z}_1; \mathbf{Z}_i)$ between the repre-

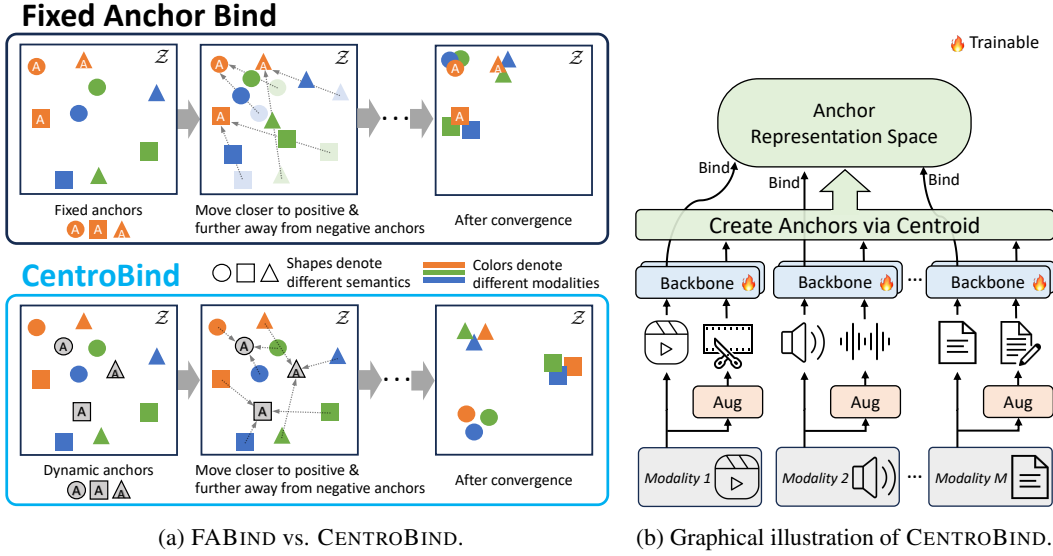


Figure 1: Fixed anchor bind methods (FABIND) binds representations to the fixed anchor modality, while CENTROBIND uses dynamic anchors. (a) Colors and shapes represent different modalities and semantic information, respectively. Z denotes the unified representation space. (b) CENTROBIND forms dynamic anchors from the centroids of positive augmentation pairs.

sensation Z_1 of the anchor modality and the representations $Z_i, i \in \{2, \dots, M\}$ of other modalities. Although these approaches are practically useful and widely adopted in learning unified multimodal representation, we theoretically and empirically demonstrate the limitations in FABIND.

Issues with fixed anchor binding. First, selecting an anchor modality is crucial but challenging, as it depends on both embedding quality and task suitability. Common choices like images or text may still be suboptimal, especially for less common modalities lacking high-quality embeddings. Second, fixing an anchor can result in lost semantic information that might be better represented by other modalities. For instance, while *text* may describe ‘a dog barks loudly,’ *sound* could reveal mood, and an *image* could add facial expression, which fixed alignment might miss. Third, optimizing only for anchor-to-other-modality overlooks information similarity between non-anchored modalities, leading to a loss of complementary insights. Addressing this among these non-anchored modalities would disrupt overall multimodal alignment, contradicting the primary goal of CENTROBIND. We formally analyze deficiencies of the fixed anchor bind approaches in Section 2.

Dynamic anchor alignment. We propose a simple yet effective modification by replacing fixed anchors with “dynamic” centroid-based anchors computed from paired samples. Our method, CENTROBIND, described in Section 3 removes the need for selecting a fixed anchor modality, instead calculates the centroid of all modality representations and generates an anchor representation, as shown in Figure 1a. Encoders are then trained to minimize the ensemble of InfoNCE loss (Oord et al., 2018) between this dynamic anchor and other representations, using the centroid as the anchor. The main intuition is that *a desirable anchor should be representative of all modalities*, capturing the most comprehensive information, with well-trained encoders producing representations that naturally cluster around this shared centroid, reflecting their underlying semantic alignment.

Our theoretical analysis demonstrates that CENTROBIND effectively addresses three critical components of multimodal learning: 1) intra-modal mutual information, 2) inter-modal mutual information, and 3) multimodal alignment via embedding similarity. By incorporating these elements within the ‘anchor alignment’ framework, CENTROBIND benefits from the significant advantages of dynamic anchors across both synthetic and real-world datasets in retrieval and classification tasks. In a sense, our approach yields an ideal unified representation space, supported by the perspective of Huh et al. (2024), which conjectures that multimodal representations align as they move toward a platonic representation—an ideal form that captures the semantic information of all modalities simultaneously.

2 PROBLEM FORMULATION

In this section, we describe general representation learning and representation binding problems in multimodal learning. Then, we analyze fixed-anchor-bind (FABIND) methods such as Image-Bind (Girdhar et al., 2023), that bind multimodal representation to a fixed modality of choice.

2.1 REPRESENTATION LEARNING FRAMEWORK

Notation. Boldface upper case letters (e.g., \mathbf{X}) denote random vectors, and a realization is denoted by the boldface lower case letters (e.g., \mathbf{x}); For $n \in \mathbb{N}$, $[n] := \{1, 2, \dots, n\}$; $P_{\mathbf{X}}$ and $P_{\mathbf{X}, \mathbf{Y}}$ denote the marginal and the joint distributions of \mathbf{X} and (\mathbf{X}, \mathbf{Y}) , respectively.

Given M datasets $\mathcal{D} = \{\mathcal{D}_i\}_{i=1}^M$, let $\mathcal{D}_i = \{(\mathbf{x}_{i,j}, \mathbf{y}_{i,j})\}_{j=1}^{N_i}$ be the dataset from the i -th modality, where $\mathbf{x}_{i,j} \in \mathcal{X}_i$ and $\mathbf{y}_{i,j} \in \mathcal{Y}_i$ are respectively the j -th input instance (e.g., feature vector) and the corresponding label in i -th modality, and we assume that $(\mathbf{x}_{i,j}, \mathbf{y}_{i,j}) \stackrel{\text{i.i.d.}}{\sim} P_{\mathbf{X}_i, \mathbf{Y}_i}$.¹ We assume that j indexes paired samples among modalities. For instance, $\mathbf{x}_{1,c}$ and $\mathbf{x}_{2,c}$ are features having similar semantic information (e.g., dog image and dog sound) in \mathcal{D}_1 and \mathcal{D}_2 . The goal of representation learning is to build M encoders $f_i : \mathcal{X}_i \rightarrow \mathcal{Z}_i$ for each modality, which maps the input instances $\mathbf{x}_{i,j}$ to its embedding $\mathbf{z}_{i,j} = f_i(\mathbf{x}_{i,j})$, preserving as much information about $\mathbf{x}_{i,j}$ as possible.

For the uni-modal case ($M = 1$), keeping maximum information about $\mathbf{x}_{1,j}$ at its embedding $\mathbf{z}_{1,j}$ is generally preferred based on the ‘‘InfoMax’’ principle (Linsker, 1988), under which the objective is to maximize mutual information $I(\mathbf{X}_i; f(\mathbf{X}_i))$ between \mathbf{X}_i and $f(\mathbf{X}_i)$. Throughout the paper, we call $I(\mathbf{X}_i; f(\mathbf{X}_i))$ *intra information* on \mathbf{X}_i . For the multimodal case ($M \geq 2$), on top of the InfoMax principle, ‘‘minimal sufficiency’’ is proposed in (Tian et al., 2020), which suggests maximizing *shared information* $I(f_i(\mathbf{X}_i); f_l(\mathbf{X}_l))$ between $f_i(\mathbf{X}_i)$ and $f_l(\mathbf{X}_l)$, while minimizing the *unique information* $I(\mathbf{X}_i; f_i(\mathbf{X}_i) | \{\mathbf{X}_l\}_{l \neq i})$. Although minimal sufficiency often leads to an efficient encoder with a better performance in numerous multimodal downstream tasks, it is not always a good strategy as there exist exceptions where the unique information on an individual modality is crucial (Liang et al., 2024b; Wang et al., 2022). In other words, the optimality of minimal sufficiency is task-dependent. To avoid task dependency, we do not consider minimal sufficiency; instead, we maximize intra and shared information without reducing unique information. Next, we formalize the notion of sufficient embedding.

Definition 1 (\mathcal{Z}_i -Sufficient embedding of \mathbf{X}_i for \mathbf{X}_l). For an embedding space \mathcal{Z}_i , the embedding $f_i(\mathbf{X}_i)$ is \mathcal{Z}_i -sufficient for \mathbf{X}_l if and only if

$$f_i \in \arg \max_{f: \mathcal{X}_i \rightarrow \mathcal{Z}_i} I(f(\mathbf{X}_i); \mathbf{X}_l), \quad (1)$$

and we call f_i sufficient encoder of \mathbf{X}_i for \mathbf{X}_l .

We note that if $i = l$, the sufficient encoder provides embeddings with maximum intra information, and if $i \neq l$, it gives embeddings with maximum shared information between i -th and l -th modalities.²

In the context of contrastive representation learning, with a goal of attaining sufficient encoders in Definition 1, InfoNCE loss $I_{\text{NCE}}(\mathbf{X}; \mathbf{Y})$ is often employed since it relates to mutual information. Specifically, InfoNCE provides a lower bound on mutual information, i.e., $I(\mathbf{X}; \mathbf{Y}) \geq -I_{\text{NCE}}(\mathbf{X}; \mathbf{Y})$ (Oord et al., 2018), thus minimizing InfoNCE leads to an increase in mutual information. InfoNCE loss between embeddings \mathbf{U} and \mathbf{V} can be written as follows:

$$I_{\text{NCE}}(\mathbf{U}; \mathbf{V} | \tau) = \mathbb{E}_{P_{\mathbf{U}, \mathbf{V}}, \prod_{i=1}^N P_{\mathbf{V}_i}} \left[-\log \frac{\exp(\mathbf{U}^\top \mathbf{V} / \tau)}{\exp(\mathbf{U}^\top \mathbf{V} / \tau) + \sum_{i=1}^N \exp(\mathbf{U}^\top \mathbf{V}_i / \tau)} \right], \quad (2)$$

where the expectation is taken with respect to $P_{\mathbf{U}, \mathbf{V}} \prod_{i=1}^N P_{\mathbf{V}_i}$. Here, we say $(\mathbf{U}, \mathbf{V}) \sim P_{\mathbf{U}, \mathbf{V}}$ a positive pair and $(\mathbf{U}, \mathbf{V}_i) \sim P_{\mathbf{U}} P_{\mathbf{V}_i}$ a negative pair. Moreover, $N \geq 1$ and $\tau > 0$ are hyper-parameters, specifying the number of negative samples and the temperature parameter. For simplicity,

¹In self-supervised learning, labels might not exist, which corresponds to the case that $\mathbf{y}_{i,j}$ are empty.

²With a proper choice of \mathcal{Z}_i ensuring $\max_{f: \mathcal{X}_i \rightarrow \mathcal{Z}_i} I(f(\mathbf{X}_i); \mathbf{X}_l) = I(\mathbf{X}_i; \mathbf{X}_l)$, Definition 1 says that $\mathbf{z}_{i,j} = f_i(\mathbf{x}_{i,j})$ is a sufficient statistic (Polyanskiy & Wu, 2024) of $\mathbf{x}_{i,j}$ for $\mathbf{x}_{l,j}$ as the encoding entails no information loss.

we assume that embeddings are normalized (Wang & Isola, 2020) and are of the same dimensionality in this paper. Under the assumption of $\tau = 1$, the exponent $\mathbf{U}^\top \mathbf{V} / \tau$ in (2) is the cosine similarity score between \mathbf{U} and \mathbf{V} .

2.2 BINDING REPRESENTATION SPACES

In addition to the objective of intra and shared information, multimodal learning often takes into account multimodal alignment (Radford et al., 2021; Duan et al., 2022). Without multimodal alignment, each modality has its own embedding structure depending on its encoder. For example, embeddings of cat and dog images, respectively, locate around $(1, 0)$ and $(0, 2)$ in \mathbb{R}^2 , whereas embeddings of cat and dog text can lie around $(0, 2)$ and $(1, 0)$. Such a misalignment can happen even for sufficient encoders (Definition 1), since the mutual information is invariant to one-to-one mappings (Polyanskiy & Wu, 2024).

To align multimodal embedding spaces, a unified representation space (Radford et al., 2021; Zhou et al., 2023) or multimodal alignment (Wang et al., 2023; Liang et al., 2024c) have gained attention and become important property in multimodal representation learning, in which embeddings of multimodal features having similar semantic should locate nearby. Among various techniques, efficient FABIND methods have been proposed (see Appendix A for a summary of FABIND methods), such as ImageBind (Girdhar et al., 2023). Their main idea is to set the image modality as a fixed anchor modality, and they minimize the InfoNCE loss between the embeddings of the anchor modality and the other modalities. More broadly, FABIND (e.g., ImageBind) aims to find encoders f_i^{FB} for all modalities except the anchor modality such that

$$f_i^{\text{FB}} = \arg \min_{f_i: \mathcal{X}_i \rightarrow \mathcal{Z}_i} I_{\text{NCE}}(f_1(\mathbf{X}_1); f_i(\mathbf{X}_i)), \forall i \in \{2, \dots, M\}, \quad (3)$$

where f_1 is the encoder for the anchor modality (e.g., image encoder in ImageBind). Note that FABIND freezes f_1 , which is initialized by an existing pretrained model, during the optimization.

2.3 ANALYSIS OF FABIND

In this section, we characterize limitations of FABIND regarding the important objective in multimodal representation learning, such as intra and shared information. To this end, we consider (3) as

$$f_i^{\text{FB}} = \arg \max_{f_i: \mathcal{X}_i \rightarrow \mathcal{Z}_i} I(f_1(\mathbf{X}_1); f_i(\mathbf{X}_i)), \forall i \in \{2, \dots, M\}, \quad (4)$$

reflecting the fact that minimizing InfoNCE loss leads to maximizing mutual information.³ Let FABIND encoders from (4) for each modality be $\mathcal{F}^{\text{FB}} = \{f_1, f_2^{\text{FB}}, \dots, f_M^{\text{FB}}\}$. The anchor encoder f_1 is fixed during the entire FABIND procedure. Moreover, we assume that $I(f_1(\mathbf{X}_1); f_i^{\text{FB}}(\mathbf{X}_i)) = I(f_1(\mathbf{X}_1); \mathbf{X}_i)$ is the maximum value that can be achieved by (4) due to data processing inequality (Polyanskiy & Wu, 2024). We next demonstrate that the quality of anchor embedding $f_1(\mathbf{X}_1)$ significantly impacts the performance of \mathcal{F}^{FB} in terms of shared information. The following propositions show the dependency of FABIND on anchor embedding quality.

Proposition 1 (FABIND with sufficient anchor). *Let $f_1^{\text{suf}}(\mathbf{X}_1)$ be a sufficient embedding of the anchor \mathbf{X}_1 , and let $\mathbf{X}_i, i \in [M]$ be discrete. Assume that $f_i^{\text{FB}}, i \in \{2, \dots, M\}$ are obtained by (4) with a sufficient anchor encoder $f_1 = f_1^{\text{suf}}$, i.e., $I(f_1^{\text{suf}}(\mathbf{X}_1); f_i^{\text{FB}}(\mathbf{X}_i)) = I(f_1^{\text{suf}}(\mathbf{X}_1); \mathbf{X}_i)$. Then,*

$$I(f_1^{\text{suf}}(\mathbf{X}_1); f_i^{\text{FB}}(\mathbf{X}_i)) = I(\mathbf{X}_1; \mathbf{X}_i), \forall i \in \{2, \dots, M\}. \quad (5)$$

Proof. The proof can be found in Appendix B.1. \square

Proposition 2 (FABIND with insufficient anchor). *Let $f_1^{\text{ins}}(\mathbf{X}_1)$ be insufficient embedding of the anchor \mathbf{X}_1 for \mathbf{X}_1 , in the sense that there exists some $\epsilon > 0$ such that $I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_1) < \epsilon \leq \max_f I(f(\mathbf{X}_1); \mathbf{X}_1)$. Assume that $f_i^{\text{FB}}, i \in \{2, \dots, M\}$ are obtained by (4) with $f_1 = f_1^{\text{ins}}$, i.e., $I(f_1^{\text{ins}}(\mathbf{X}_1); f_i^{\text{FB}}(\mathbf{X}_i)) = I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_i)$. Then,*

$$I(f_1^{\text{ins}}(\mathbf{X}_1); f_i^{\text{FB}}(\mathbf{X}_i)) < \epsilon, \forall i \in \{2, \dots, M\}. \quad (6)$$

³In contrast to (3), f_i^{FB} in (4) might not be aligned with other modalities due to the one-to-one mapping invariant property of mutual information. However, we here do not analyze multimodal alignment of FABIND from (4), but rather investigate the performance of encoders in terms of the sufficiency in Definition 1.

Proof. The proof can be found in Appendix B.2. \square

Proposition 1 shows that the FABIND encoders \mathcal{F}^{FB} learned with a sufficient anchor embedding can achieve the maximum shared information between the anchor and the other modalities. However, it does not guarantee the shared information between *non-anchored* modalities $I(f_i(\mathbf{X}_i); f_l(\mathbf{X}_l))$, $i, l \neq 1$, which can also be seen from (4). Proposition 2 establishes that an insufficient anchor may lead to a reduction of shared information between the anchor and the other modalities, implying that the performance of FABIND solely depends on the quality of the anchor.

The analysis reveals several limitations in FABIND. Firstly, achieving maximum shared information requires sufficient anchor representation, which depends on having both an informative modality and a sufficient encoder. Without these conditions, FABIND may not effectively capture shared information. Secondly, even with sufficient anchor representation, FABIND may not provide encoders with maximum intra information. This is because its objective function (4) does not take into account intra information. Thirdly, the objective function of FABIND (4) focuses solely on learning shared information between pairs of anchor and non-anchored modalities, while disregarding shared information among non-anchor modalities. *It implies that FABIND does not necessitate to capture shared information among non-anchored modalities.* This limitation renders FABIND less effective when significant shared information exists among non-anchored modalities. Lastly, the representation produced by FABIND may not approximate an ideal representation, such as ‘‘Platonic representation’’ (Huh et al., 2024). The Platonic representation is an idealized depiction of reality, which generates all modalities through projections. While integrating all modalities is crucial for constructing such a comprehensive representation, FABIND approach, which focuses solely on the anchor modality, falls short of this ideal. A more promising direction, which we pursue in this paper, involves learning representations from multiple modalities, fully leveraging fine-grained, sample-level information.

In summary, FABIND exhibits the following weaknesses:

- P1:** over-reliance on a single anchor modality;
- P2:** failure to capture intra information;
- P3:** absence of shared information among non-anchored modalities;
- P4:** fundamental limitation of a fixed anchor in representativeness of all modalities.

Next, we propose our method CENTROBIND, which does not require any choice of anchor modality and is able to capture intra and shared information.

3 TOWARD A DESIRABLE UNIFIED REPRESENTATION SPACE

Our method, CENTROBIND, is a variant of Bind methods, but does not require choosing a single anchor modality. The main intuition deriving CENTROBIND is as follows: 1) A desirable multimodal embedding should capture representation of all modalities; 2) A desirable unified embedding should attain the highest average alignment. To this end, we generate the anchor representation from all modalities by calculating the centroid of the representations of the modalities. Then, we train the encoders toward minimizing the InfoNCE loss between anchor and other modalities, which is identical to FABIND with the generated anchor instead of a fixed anchor. Here we formally describe CENTROBIND, and show that the method aligns multimodal representations and simultaneously maximizes intra and shared information.

3.1 CENTROBIND

To present a general framework for CENTROBIND, we consider M modalities with corresponding encoders $\{f_i\}_{i=1}^M$. The algorithmic presentation of CENTROBIND is in Algorithm 1, and a graphical illustration is given in Figure 1b. In the following, we elaborate detailed steps of CENTROBIND.

Initial encoders. We initialize M encoders $f_i : \mathcal{X}_i \rightarrow \mathcal{Z}$, $\forall i \in [M]$ for the M modalities. These encoders can either be pretrained models (i.e., backbones) or parameterized models with random weights. The primary constraint for these initial encoders is that their output space must be \mathcal{Z} . This

Algorithm 1 CENTROBIND

```

1: Initialize encoders  $f_1^{(0)}, f_2^{(0)}, \dots, f_M^{(0)}$ .
2: for  $t = 0, 1, \dots, t_{\max}$  do
3:   Sample a batch dataset  $B$  from multimodal datasets  $\{\mathcal{D}_i\}_i$ .
4:   Generate anchor embeddings  $\{\mathbf{a}_j\}_{j \in \mathcal{I}_B}$  using  $\mathbf{a}_j = \text{mean}(\{f_i^{(t)}(\mathbf{x}'_{i,j})\}_i)$  in (7)
5:   for  $i = 1, \dots, M$  do
6:     Optimize  $f_i^{(t+1)}$  toward minimizing  $\mathcal{L}_{\text{CB}}(f_i^{(t+1)}|\tau)$  in (8)
7:   end for
8: end for

```

ensures that all data samples are mapped to the same output space \mathcal{Z} . When using pretrained encoders that produce embeddings in different output spaces, we attach projection layers to standardize the shapes of the embeddings, ensuring consistency of output space across modalities.

Anchor embedding. Recall that $\mathbf{x}_{i,j} \in \mathcal{D}_i$ denotes the j -th feature in the i -th modality, where j indexes positive pairs of features (e.g., different views of the same object). In each training iteration of CENTROBIND, we need to compute an anchor embedding \mathbf{a}_j for the j -th multimodal positive features $\{\mathbf{x}_{i,j}\}_{i=1}^M$. This anchor \mathbf{a}_j serves as a desirable aligned embedding for these features. The anchor \mathbf{a}_j is calculated as follows:

$$\mathbf{a}_j = \text{mean}(\{f_i(\mathbf{x}'_{i,j})\}_i), \quad (7)$$

where $\text{mean}(\cdot)$ denotes the mean operator that computes the average of its input, and $\mathbf{x}'_{i,j}$ represents an augmented version of $\mathbf{x}_{i,j}$. If $\{\mathbf{x}_{i,j}\}_{i=1}^M$ are available in multimodal datasets, the anchor is given by $\mathbf{a}_j = \frac{1}{M} \sum_{i=1}^M f_i(\mathbf{x}'_{i,j})$. If only $m < M$ positive pairs are present among M modalities, the anchor is given by $\mathbf{a}_j = \frac{1}{m} \sum_{i \in \mathcal{I}_j} f_i(\mathbf{x}'_{i,j})$, where \mathcal{I}_j is the set of indices of modalities having the m available features.

Binding encoders to the anchor. Once anchor embeddings $\{\mathbf{a}_j\}_j$ are derived from a batch of data $B = \{\mathbf{x}_{i,j}\}_{i,j}$, CENTROBIND aligns each modality-specific encoder embedding with the anchor embedding by minimizing the InfoNCE loss. Specifically, let $\mathbf{A} = \text{mean}(\{f_i(\mathbf{X}_i)\}_i)$ represent the anchor embedding variable. Then, CENTROBIND aims to minimize the InfoNCE loss $I_{\text{NCE}}(\mathbf{A}; f_i(\mathbf{X}_i))$ across all modalities $i \in [M]$. A detailed expression for this loss is provided in (8).

For symmetry, CENTROBIND optimizes the following loss function:

$$\mathcal{L}_{\text{CB}}(f_i|\tau) = I_{\text{NCE}}(\mathbf{A}; f_i(\mathbf{X}_i)|\tau) + I_{\text{NCE}}(f_i(\mathbf{X}_i); \mathbf{A}|\tau), \quad (8)$$

where $\mathcal{L}_{\text{CB}}(f_i|\tau)$ denotes the loss function for the i -th modality. In particular, with a batch data $B = \{\mathbf{x}_{i,j} : i \in [M], j \in \mathcal{I}_B\}$, the loss can be computed as

$$I_{\text{NCE}}(\mathbf{A}; f_i(\mathbf{X}_i)|\tau) = -\frac{1}{|\mathcal{I}_B|} \sum_{k=1}^{|\mathcal{I}_B|} \log \frac{\exp(\mathbf{a}_k^\top f_i(\mathbf{x}_{i,k})/\tau)}{\sum_{j \in \mathcal{I}_B} \exp(\mathbf{a}_k^\top f_i(\mathbf{x}_{i,j})/\tau)} \quad \text{and} \quad (9a)$$

$$I_{\text{NCE}}(f_i(\mathbf{X}_i); \mathbf{A}|\tau) = -\frac{1}{|\mathcal{I}_B|} \sum_{k=1}^{|\mathcal{I}_B|} \log \frac{\exp(\mathbf{a}_k^\top f_i(\mathbf{x}_{i,k})/\tau)}{\sum_{j \in \mathcal{I}_B} \exp(f_i^\top(\mathbf{x}_{i,k})\mathbf{a}_j/\tau)}. \quad (9b)$$

3.2 THEORETICAL ANALYSIS OF CENTROBIND

We start by providing a lower bound on the objective function of CENTROBIND $\mathcal{L}_{\text{CB}}(f_i|\tau)$ (8) in Theorem 1, followed by an analysis of minimizing $\mathcal{L}_{\text{CB}}(f_i|\tau)$.

Theorem 1. Consider $B = \{\mathbf{x}_{i,j} : i \in [M], j \in \mathcal{I}_B\}$ with a set of indices \mathcal{I}_B , where $\mathbf{x}_{i,j}$ is the j -th sample of i -th modality. Then, for any encoders $\{f_i\}_i$ and for any $\tau > 0$, (9a) is bounded as

$$I_{\text{NCE}}(\mathbf{A}; f_i(\mathbf{X}_i) | \tau) \geq \frac{1}{|\mathcal{I}_B|} \sum_{l=1}^M I_{\text{NCE}} \left(f_l(\mathbf{X}'_l); f_i(\mathbf{X}_i) \mid \frac{\tau M}{|\mathcal{I}_B|} \right) - \frac{1}{|\mathcal{I}_B|} \sum_{k=1}^{|\mathcal{I}_B|} \log C_{\mathcal{F},k,i}, \quad (10)$$

where $C_{\mathcal{F},k,i} = \frac{(c_{\mathcal{F},k,i}^{\min} + c_{\mathcal{F},k,i}^{\max})^2}{4c_{\mathcal{F},k,i}^{\min}c_{\mathcal{F},k,i}^{\max}}$ with $g(l,j|k,i) := \exp\left(\frac{|\mathcal{I}_B|f_l^\top(\mathbf{x}'_{l,k})f_i(\mathbf{x}_{i,j})}{\tau M}\right)$,

$$c_{\mathcal{F},k,i}^{\min} = \min_{l \in [M], j \in \mathcal{I}_B} g(l,j|k,i), \text{ and } c_{\mathcal{F},k,i}^{\max} = \max_{l \in [M], j \in \mathcal{I}_B} g(l,j|k,i). \quad (11)$$

Proof. The proof is provided in Appendix B.3. \square

Theorem 1 provides a lower bound of $I_{\text{NCE}}(\mathbf{A}; f_i(\mathbf{X}_i) \mid \tau)$ in (9a), which is a part of CENTROBIND objective $\mathcal{L}_{\text{CB}}(f_i|\tau)$. Thus CENTROBIND consequentially minimizes the lower bound (10) that consists of two terms, $\sum_{l=1}^M I_{\text{NCE}}\left(f_l(\mathbf{X}'_l); f_i(\mathbf{X}_i) \mid \frac{\tau M}{|\mathcal{I}_B|}\right)$ and $-\sum_{k=1}^{|\mathcal{I}_B|} \log C_{\mathcal{F},k,i}$. We next provide some intuitive explanation on the effect of such a minimization of the lower bound.

The effect of minimizing $\sum_{l=1}^M I_{\text{NCE}}\left(f_l(\mathbf{X}'_l); f_i(\mathbf{X}_i) \mid \frac{\tau M}{|\mathcal{I}_B|}\right)$. The objective of minimizing $\sum_{l=1}^M I_{\text{NCE}}(f_l(\mathbf{X}'_l); f_i(\mathbf{X}_i) \mid \frac{\tau M}{|\mathcal{I}_B|})$ involves reducing several InfoNCE losses. Here, each term in the sum represents the InfoNCE loss between embeddings $f_l(\mathbf{X}'_l)$ from modality l and $f_i(\mathbf{X}_i)$ from modality i , with $\frac{\tau M}{|\mathcal{I}_B|}$ being the temperature parameter for scaling the loss. This summation can be divided into two components: 1) Intra Information: When $l = i$, the term measures the similarity between embeddings within the same modality. Minimizing this loss enhances the representation of modality i , improving intra information; 2) Shared Information: When $l \neq i$, the term measures the similarity between embeddings from different modalities. Minimizing these losses helps in learning shared information between modalities, contributing to a more coherent multimodal representation.

By optimizing this summation, CENTROBIND effectively captures both intra and shared information. This results in a more accurate representation for each modality. In contrast, as noted in Section 2.3, FABIND does not adequately capture intra information and shared information between non-anchored modalities. This limitation highlights the advantage of CENTROBIND in achieving a more integrated multimodal representation than fixed anchor binding methods.

The effect of minimizing $-\sum_{k=1}^{|\mathcal{I}_B|} \log C_{\mathcal{F},k,i}$. We show the effect of growing $C_{\mathcal{F},k,i}$ in terms of cosine similarity score between embeddings. Since $C_{\mathcal{F},k,i} = \frac{1}{4} \left(\sqrt{\gamma} + \sqrt{\frac{1}{\gamma}} \right)^2$ with $\gamma = \frac{c_{\mathcal{F},k,i}^{\max}}{c_{\mathcal{F},k,i}^{\min}} \geq 1$, maximizing $C_{\mathcal{F},k,i}$ is equivalent to simultaneously maximizing $c_{\mathcal{F},k,i}^{\max}$ and minimizing $c_{\mathcal{F},k,i}^{\min}$. For ease of the analysis, we assume that the encoders are reasonably well-trained. Then, since a positive pair of embeddings normally yields higher similarity score, $c_{\mathcal{F},k,i}^{\max}$ should be obtained by choosing $l = i$ and $j = k$ in (11) as such choices make $\mathbf{x}'_{l,k}$ be positive pair with $\mathbf{x}_{i,j}$. Thus, $c_{\mathcal{F},k,i}^{\max}$ is roughly proportional to the similarity score of a positive pair of embeddings. Conversely, $c_{\mathcal{F},k,i}^{\min}$ corresponds to the similarity scores of negative pairs, which tend to be low. Hence, minimizing $-\sum_{k=1}^{|\mathcal{I}_B|} \log C_{\mathcal{F},k,i}$ enhances the similarity scores for positive pairs and reduces those for negative pairs, improving the overall multimodal alignment.

The preceding analyses demonstrate that CENTROBIND addresses the limitations **P1**, **P2**, **P3**, and **P4** of FABIND identified in Section 2.3. We now conclude the analysis with a discussion that the unified representation of CENTROBIND is likely closer to an ideal platonic representation (Huh et al., 2024) compared to FABIND's. A platonic representation is defined as an ideal representation of reality that induces information in all modalities. From this perspective, a representation derived solely from a single modality, without leveraging others, may not be sufficient. In contrast to FABIND's approach, which relies exclusively on the fixed anchor modality, CENTROBIND constructs its unified representation space by incorporating information from all modalities. This suggests that CENTROBIND's unified space is likely to retain a more comprehensive representation of all modalities. Thus, CENTROBIND emerges as a promising approach for developing a Planotic representation.

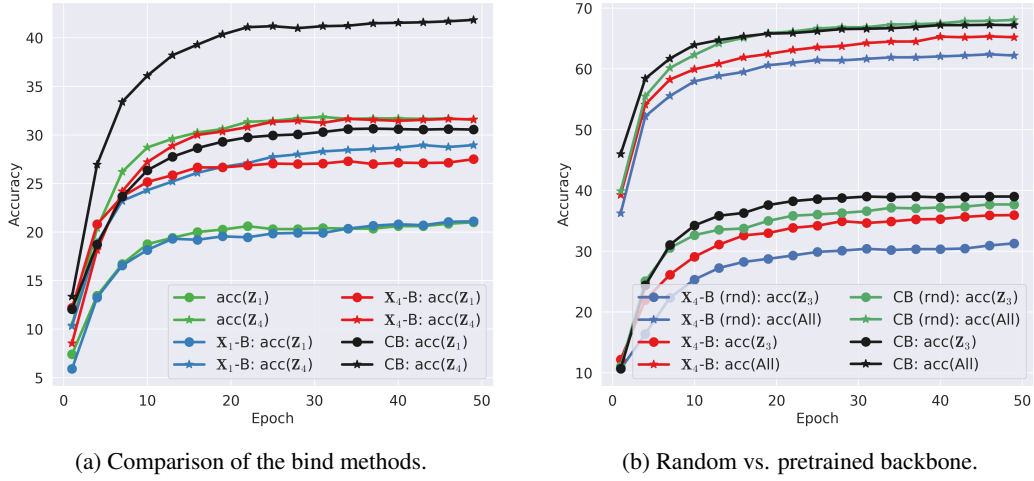


Figure 2: Accuracy as a measure of the representation space quality. Abbreviation: X_i-B or CB: applying FABIND with anchor X_i or applying CentroBind; acc(Z_i) or acc(All): accuracy of Z_i or of concatenated embeddings (Z₁, ..., Z_M); (rnd): if random backbones are used.

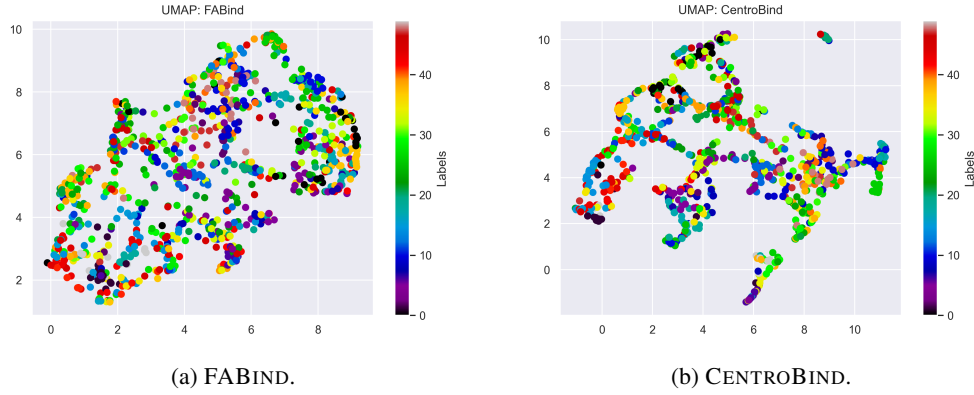


Figure 3: Representation visualization via UMAP.

4 EXPERIMENT

4.1 EXPERIMENTS WITH SYNTHETIC DATASET

Synthetic datasets. We employ a latent variable model (Bishop & Nasrabadi, 2006) for generating synthetic multimodal datasets. A latent variable model is a statistical model for data $\mathbf{X} \in \mathbb{R}^{d_x}$, under which \mathbf{X} is generated according to a conditional probability distribution $P_{\mathbf{X}|\mathbf{Z}}$, where $\mathbf{Z} \in \mathbb{R}^{d_z}$ is the latent variable. In terms of the representation learning framework, \mathbf{Z} can be seen as the true representation of \mathbf{X} . We assume that the class label $\mathbf{Y} \in [K]$ and the latent variable \mathbf{Z} are jointly distributed according to $P_{\mathbf{Z}, \mathbf{Y}}$. In our setting, we exploit Gaussian mixture model (GMM) (Bishop & Nasrabadi, 2006) for the latent variable \mathbf{Z} , and we generate M modalities $\mathbf{X}_i = g_i(\mathbf{Z}) + \mathbf{N}$, $i \in [M]$ with random noise \mathbf{N} and some non-linear projections $g_i : \mathbb{R}^{d_z} \rightarrow \mathbb{R}^{d_x}$. We choose the projections in a way such that each model can be ranked in ascending order, i.e., \mathbf{X}_1 is the worst, and \mathbf{X}_4 is the best modality in terms of their inherent correlation with the latent variable. The class label \mathbf{Y} is set to the component id of GMM (for details, see Appendix C.1).

Experiment results. Figure 2 shows the classification accuracies with a synthetic dataset of $M = 4$ modalities. To obtain the results in Figure 2a, we initialize pretrained backbones for all modalities, apply FABIND (X_i-B) with anchor X_i or CENTROBIND (CB), and evaluate accuracy (acc(Z_i)) with embeddings from i -th modality. We provide acc(Z_i), without any binding, for a reference.

Figure 2a verifies our analysis of FABIND (Section 2.3) and CentroBind (Section 3.2): (1) the comparison between \mathbf{X}_1 -B and \mathbf{X}_4 -B shows the importance of choosing anchor modality; (2) the comparison between $\text{acc}(\mathbf{Z}_4)$ and \mathbf{X}_1 -B: $\text{acc}(\mathbf{Z}_4)$ shows a performance deterioration by FABIND, demonstrating that FABIND does not capture intra information; (3) CB consistently outperform all FABIND, indicating that CB successfully captures elements that FABIND overlooks, including intra information and shared information among non-anchored modalities.

Figure 2b includes accuracies of FABIND and CB with random backbones. Similarly, CB outperforms all baselines. Somewhat surprisingly, CB with random backbones (green curves) also performs better than FABIND with pretrained backbones (red curves). This further supports our analysis that CENTROBIND is robust to the backbone quality as it optimizes intra and shared information, whereas FABIND is sensitive to the backbone quality. Overall, these empirical results validate our findings. We provide additional experimental results on synthetic datasets with $M = 6, 8$ in Appendix C.1. With the larger number of modalities, CB still outperforms the baselines, strengthening our analysis of CB.

In addition, we visualize the embeddings learned by FABIND and CENTROBIND using UMAP (McInnes et al., 2018) in Figure 3 (a visualization using t-SNE (Van der Maaten & Hinton, 2008) is in Figure 5 in Appendix C.1). For this visualization, we generate the synthetic datasets with 4 modalities such that each modality is equally informative, and plot the embeddings for \mathbf{X}_1 . FABIND is anchored at \mathbf{X}_4 and both binding methods use pre-trained backbones. Figure 3 shows that CENTROBIND embeddings are better clustered, while FABIND embeddings are scattered, implying that CENTROBIND constructs better embedding structure than FABIND.

Although convergence rate is not a critical aspect for representation learning which is usually done offline, we empirically explore convergence speed for both CENTROBIND and FABIND. In Figure 4, we plot train loss curves, and it shows that the loss curves of both methods behave similarly, implying that dynamic anchor does not pose a problem regarding convergence or stability. We discuss the detail in Appendix C.1.

Lastly, instead of centroid we compare other possible dynamic anchor generation methods, such as weighted average and random anchor in Appendix C.1. Intuitively, if modality imbalance problem exists, weighted average or other methods may lead to better dynamic anchor generation since it leverages additional information on downstream tasks or datasets. Nevertheless, CENTROBIND performs well even if it does not necessitate addition extra information. In summary, CENTROBIND generally performs better or similar to the weighted average dynamic anchor generation. As our goal is multimodal alignment in representation learning without domain knowledge and downstream tasks, centroid is a reasonable dynamic anchor generation.

4.2 EXPERIMENTS WITH REAL-WORLD DATASET

In this section, we provide experiment results with a real-world dataset. We compare CENTROBIND and FABIND anchored at text modality. We utilize the MUsTARD dataset (Castro et al., 2019) for its rich combination of multimodal data with more than two modalities. It consists of 690 video clips (including audio) and text for sarcasm detection with labels such as sarcasm indicators and speaker names. For the backbones in FABIND and CENTROBIND, we use the pretrained VideoMAE model (Tong et al., 2022) for video data, the pretrained WaveLM model (Chen et al., 2022) for audio data, and the pretrained BERT model (Devlin et al., 2019) for text data. A detailed description of the training setting is provided in Appendix C.2.

Downstream tasks. We perform evaluations in zero-shot binary and multi-class classification tasks, One-to-One cross-modal retrieval, and Two-to-One cross-modal retrieval. For classification tasks, we use a Multi-Layer Perceptron (MLP) to perform sarcasm detection as a binary classification and speaker classification with 23 multi-class categories. In particular, MLP is trained on embeddings in a single modality (denoted by \mathbf{Tr} in Table 2) and accuracy is evaluated on another modality (denoted by \mathbf{Ev} in Table 2). In retrieval tasks, we measure the accuracy of correct retrieval. For One-to-One case, we retrieve data sample in different modality by choosing the closest embedding from a single input embedding, while for Two-to-One case we choose the closest embedding from the centroid of two input embeddings in two modalities. We denote input and target modalities with \rightarrow in Table 1.

Table 1: Zero-shot one-to-one and two-to-one retrieval accuracy. (\mathcal{V} : video, \mathcal{A} : audio, \mathcal{T} : text)

One-to-One					Two-to-One				
Method	Retrieval	Top-1	Top-5	Top-10	Method	Retrieval	Top-1	Top-5	Top-10
FABIND	$\mathcal{V} \rightarrow \mathcal{T}$	0.446	0.719	0.822	FABIND	$\mathcal{V} \mathcal{A} \rightarrow \mathcal{T}$	0.309	0.665	0.781
CENTROBIND	$\mathcal{V} \rightarrow \mathcal{T}$	0.483	0.764	0.850	CENTROBIND	$\mathcal{V} \mathcal{A} \rightarrow \mathcal{T}$	0.745	0.957	0.978
FABIND	$\mathcal{A} \rightarrow \mathcal{T}$	0.077	0.238	0.367	FABIND	$\mathcal{T} \mathcal{A} \rightarrow \mathcal{V}$	0.180	0.401	0.513
CENTROBIND	$\mathcal{A} \rightarrow \mathcal{T}$	0.233	0.517	0.678	CENTROBIND	$\mathcal{T} \mathcal{A} \rightarrow \mathcal{V}$	0.388	0.646	0.768
FABIND	$\mathcal{T} \rightarrow \mathcal{V}$	0.812	0.946	0.978	FABIND	$\mathcal{T}, \mathcal{V} \rightarrow \mathcal{A}$	0.099	0.257	0.364
CENTROBIND	$\mathcal{T} \rightarrow \mathcal{V}$	0.591	0.839	0.909	CENTROBIND	$\mathcal{T}, \mathcal{V} \rightarrow \mathcal{A}$	0.232	0.490	0.625
FABIND	$\mathcal{A} \rightarrow \mathcal{V}$	0.058	0.154	0.226					
CENTROBIND	$\mathcal{A} \rightarrow \mathcal{V}$	0.052	0.184	0.284					
FABIND	$\mathcal{T} \rightarrow \mathcal{A}$	0.201	0.438	0.584					
CENTROBIND	$\mathcal{T} \rightarrow \mathcal{A}$	0.290	0.572	0.706					
FABIND	$\mathcal{V} \rightarrow \mathcal{A}$	0.051	0.155	0.223					
CENTROBIND	$\mathcal{V} \rightarrow \mathcal{A}$	0.054	0.175	0.258					

Results on cross-modal retrieval. Table 1 shows the performance for one-to-one and two-to-one retrieval tasks. CENTROBIND consistently excels in one-to-one retrieval for text and audio modalities, while FABIND performs better for video retrieval. This might be due to a strong dependency between text and video, which may be suitable for FABIND anchored at text modality. A notable observation is that the centroid of video and audio embeddings achieves the best text retrieval performance including one-to-one retrieval. This implies complementary information exists and is captured by CENTROBIND.

Results on sarcasm & speaker classification.

Table 2 presents results for sarcasm detection and speaker classification tasks, where Sar-1 indicates Top-1 accuracy for sarcasm, and Spk- k , $k = 1, 3, 5$ represent Top- k accuracies for speaker classification. In this experiment, CENTROBIND consistently outperforms FABIND across all pairs of train and evaluation modalities, which can be distributed to CENTROBIND generally learning a better shared embedding space than FABIND. Although a direct comparison is not feasible, we present the sarcasm detection accuracy of FactorCL (Liang et al., 2024b) and SimMMDG (Dong et al., 2023) approaches for reference. It is important to highlight that CENTROBIND and FABIND are trained on a single modality (**Tr**) and evaluated on a different modality (**Ev**) in a zero-shot setting. In contrast, FactorCL and SimMMDG employ supervised contrastive learning, utilizing sarcasm labels during pretraining to capture task-specific features. At inference, all three modalities are simultaneously used to predict the class. As analyzed in Section 2.3 and Section 3.2, these results highlight the CENTROBIND’s ability to preserve intra and shared information among modalities, which are useful in unknown downstream tasks. Moreover, the zero-shot setting verifies the multimodal alignment of both methods.

Table 2: Accuracy results for Sarcasms and Speakers. (Tr: training modality, Ev: evaluation modality, \mathcal{V} : video, \mathcal{A} : audio, \mathcal{T} : text). Asterisks denote accuracy evaluated in different settings.

Method	Tr, (Ev)	Sar-1	Spk-1	Spk-3	Spk-5
FABIND	$\mathcal{V}, (\mathcal{T})$	0.572	0.243	0.445	0.630
CENTROBIND		0.694	0.368	0.670	0.791
FABIND	$\mathcal{A}, (\mathcal{T})$	0.535	0.241	0.509	0.635
CENTROBIND		0.655	0.346	0.610	0.741
FABIND	$\mathcal{T}, (\mathcal{V})$	0.706	0.378	0.614	0.730
CENTROBIND		0.716	0.474	0.736	0.836
FABIND	$\mathcal{A}, (\mathcal{V})$	0.604	0.255	0.472	0.636
CENTROBIND		0.626	0.326	0.548	0.703
FABIND	$\mathcal{V}, (\mathcal{A})$	0.623	0.228	0.484	0.628
CENTROBIND		0.683	0.243	0.475	0.632
FABIND	$\mathcal{T}, (\mathcal{A})$	0.648	0.186	0.455	0.577
CENTROBIND		0.691	0.290	0.546	0.714
FactorCL*	$\mathcal{V}, \mathcal{A}, \mathcal{T}$	0.699	-	-	-
SimMMDG*	$\mathcal{V}, \mathcal{A}, \mathcal{T}$	0.725	-	-	-

5 CONCLUSIONS

In this paper, we mathematically analyze the limitations of fixed-anchor-bind methods (FABIND), including over-reliance on the choice of anchor modality, and failing to capture both intra and shared information among non-anchored modalities. To overcome such shortcomings, we propose CENTROBIND, which aligns multimodal embeddings to dynamic anchors constructed by centroids of the embeddings, hence removing the need for anchor modality. Moreover, we theoretically study CENTROBIND, showing that it captures intra- and shared information. Extensive experiments on both synthetic and real-world datasets show that CENTROBIND significantly outperforms FABIND, providing a robust unified representation space and validating our analysis on CENTROBIND and FABIND.

REFERENCES

- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- Yusuf Aytar, Carl Vondrick, and Antonio Torralba. See, hear, and read: Deep aligned representations. *arXiv preprint arXiv:1706.00932*, 2017.
- Niels Balemans, Ali Anwar, Jan Steckel, and Siegfried Mercelis. Lidar-bind: Multi-modal sensor fusion through shared latent embeddings. *IEEE Robotics and Automation Letters*, 2024.
- Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. *IEEE transactions on pattern analysis and machine intelligence*, 41(2): 423–443, 2018.
- Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- Santiago Castro, Devamanyu Hazarika, Verónica Pérez-Rosas, Roger Zimmermann, Rada Mihalcea, and Soujanya Poria. Towards multimodal sarcasm detection (an _obviously_ perfect paper). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4619–4629, 2019.
- Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6): 1505–1518, 2022.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Tamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- Aayush Dhakal, Subash Khanal, Srikumar Sastry, Adeel Ahmad, and Nathan Jacobs. Geobind: Binding text, image, and audio through satellite images. *arXiv preprint arXiv:2404.11720*, 2024.
- Hao Dong, Ismail Nejjar, Han Sun, Eleni Chatzi, and Olga Fink. Simmmdg: A simple and effective framework for multi-modal domain generalization. *Advances in Neural Information Processing Systems*, 36:78674–78695, 2023.
- Chenzhuang Du, Jiaye Teng, Tingle Li, Yichen Liu, Tianyuan Yuan, Yue Wang, Yang Yuan, and Hang Zhao. On uni-modal feature learning in supervised multi-modal learning. In *International Conference on Machine Learning*, pp. 8632–8656. PMLR, 2023.
- Jiali Duan, Liqun Chen, Son Tran, Jinyu Yang, Yi Xu, Belinda Zeng, and Trishul Chilimbi. Multi-modal alignment using representation codebook. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15651–15660, 2022.
- Bruno Dumas, Jonathan Pirau, and Denis Lalanne. Modelling fusion of modalities in multimodal interactive systems with mmmm. In *Proceedings of the 19th ACM International Conference on Multimodal Interaction*, pp. 288–296, 2017.
- Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- Yuan Gao, Sangwook Kim, David E Austin, and Chris McIntosh. Medbind: Unifying language and multimodal medical data embeddings, 2024. URL <https://arxiv.org/abs/2403.12894>.

- Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15180–15190, 2023.
- Dalu Guo, Chang Xu, and Dacheng Tao. Bilinear graph networks for visual question answering. *IEEE Transactions on neural networks and learning systems*, 34(2):1023–1034, 2021.
- Ziyu Guo, Renrui Zhang, Xiangyang Zhu, Yiwen Tang, Xianzheng Ma, Jiaming Han, Kexin Chen, Peng Gao, Xianzhi Li, Hongsheng Li, et al. Point-bind & point-llm: Aligning point cloud with multi-modality for 3d understanding, generation, and instruction following. *arXiv preprint arXiv:2309.00615*, 2023.
- Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas Dengel. Audioclip: Extending clip to image, text and audio. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 976–980. IEEE, 2022.
- David Harwath, Adria Recasens, Dídac Surís, Galen Chuang, Antonio Torralba, and James Glass. Jointly discovering visual objects and spoken words from raw sensory input. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 649–665, 2018.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.
- Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. Position: The platonic representation hypothesis. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 20617–20642. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/huh24a.html>.
- Paul Pu Liang, Yun Cheng, Xiang Fan, Chun Kai Ling, Suzanne Nie, Richard Chen, Zihao Deng, Nicholas Allen, Randy Auerbach, Faisal Mahmood, et al. Quantifying & modeling multimodal interactions: An information decomposition framework. *Advances in Neural Information Processing Systems*, 36, 2024a.
- Paul Pu Liang, Zihao Deng, Martin Q Ma, James Y Zou, Louis-Philippe Morency, and Ruslan Salakhutdinov. Factorized contrastive learning: Going beyond multi-view redundancy. *Advances in Neural Information Processing Systems*, 36, 2024b.
- Paul Pu Liang, Akshay Goindani, Talha Chafekar, Leena Mathur, Haoifei Yu, Ruslan Salakhutdinov, and Louis-Philippe Morency. Hemm: Holistic evaluation of multimodal foundation models. *arXiv preprint arXiv:2407.03418*, 2024c.
- Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. Foundations & trends in multimodal machine learning: Principles, challenges, and open questions. *ACM Computing Surveys*, 56(10): 1–42, 2024d.
- R. Linsker. Self-organization in a perceptual network. *Computer*, 21(3):105–117, 1988. doi: 10.1109/2.36.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- Yuanhuiyi Lyu, Xu Zheng, Jiazhou Zhou, and Lin Wang. Unibind: Llm-augmented unified and balanced representation space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 26752–26762, June 2024.
- Divyam Madaan, Taro Makino, Sumit Chopra, and Kyunghyun Cho. A framework for multi-modal learning: Jointly modeling inter- & intra-modality dependencies. *arXiv preprint arXiv:2405.17613*, 2024.

- L. McInnes, J. Healy, and J. Melville. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *ArXiv e-prints*, February 2018.
- Arsha Nagrani, Shan Yang, Anurag Arnab, Aren Jansen, Cordelia Schmid, and Chen Sun. Attention bottlenecks for multimodal fusion. *Advances in neural information processing systems*, 34:14200–14213, 2021.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Xiaokang Peng, Yake Wei, Andong Deng, Dong Wang, and Di Hu. Balanced multimodal learning via on-the-fly gradient modulation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8238–8247, 2022.
- Yury Polyanskiy and Yihong Wu. *Information Theory: From Coding to Learning*. Cambridge University Press, 2024.
- 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 *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Daniel Seichter, Mona Köhler, Benjamin Lewandowski, Tim Wengelfeld, and Horst-Michael Gross. Efficient rgb-d semantic segmentation for indoor scene analysis. In *2021 IEEE international conference on robotics and automation (ICRA)*, pp. 13525–13531. IEEE, 2021.
- Yuki Seo. Generalized Pólya–Szegő type inequalities for some non-commutative geometric means. *Linear Algebra and its Applications*, 438(4):1711–1726, 2013.
- Sheng Shen, Chunyuan Li, Xiaowei Hu, Yujia Xie, Jianwei Yang, Pengchuan Zhang, Zhe Gan, Lijuan Wang, Lu Yuan, Ce Liu, et al. K-lite: Learning transferable visual models with external knowledge. *Advances in Neural Information Processing Systems*, 35:15558–15573, 2022.
- Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, 2019.
- Jiajie Teng, Huiyu Duan, Yucheng Zhu, Sijing Wu, and Guangtao Zhai. Mvbind: Self-supervised music recommendation for videos via embedding space binding, 2024. URL <https://arxiv.org/abs/2405.09286>.
- Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What makes for good views for contrastive learning? *Advances in neural information processing systems*, 33:6827–6839, 2020.
- Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *Advances in neural information processing systems*, 35:10078–10093, 2022.
- Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J. Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6558–6569, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1656. URL <https://aclanthology.org/P19-1656>.
- Xinming Tu, Zhi-Jie Cao, xia chenrui, Sara Mostafavi, and Ge Gao. Cross-linked unified embedding for cross-modality representation learning. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 15942–15955. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/662b1774ba8845fc1fa3d1fc0177ceeb-Paper-Conference.pdf.
- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.

- Fei Wang, Liang Ding, Jun Rao, Ye Liu, Li Shen, and Changxing Ding. Can linguistic knowledge improve multimodal alignment in vision-language pretraining? *arXiv preprint arXiv:2308.12898*, 2023.
- Haoqing Wang, Xun Guo, Zhi-Hong Deng, and Yan Lu. Rethinking minimal sufficient representation in contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16041–16050, 2022.
- Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through alignment and uniformity on the hypersphere. In *International conference on machine learning*, pp. 9929–9939. PMLR, 2020.
- Weiyao Wang, Du Tran, and Matt Feiszli. What makes training multi-modal classification networks hard? In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12695–12705, 2020a.
- Zehan Wang, Ziang Zhang, Xize Cheng, Rongjie Huang, Luping Liu, Zhenhui Ye, Haifeng Huang, Yang Zhao, Tao Jin, Peng Gao, and Zhou Zhao. Freebind: Free lunch in unified multimodal space via knowledge fusion. In *Forty-first International Conference on Machine Learning*, 2024a.
- Zehan Wang, Ziang Zhang, Hang Zhang, Luping Liu, Rongjie Huang, Xize Cheng, Hengshuang Zhao, and Zhou Zhao. Omnibind: Large-scale omni multimodal representation via binding spaces. *arXiv preprint arXiv:2407.11895*, 2024b.
- Zhecheng Wang, Haoyuan Li, and Ram Rajagopal. Urban2vec: Incorporating street view imagery and pois for multi-modal urban neighborhood embedding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 1013–1020, 2020b.
- Alex Wilf, Martin Q Ma, Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. Face-to-face contrastive learning for social intelligence question-answering. In *2023 IEEE 17th International Conference on Automatic Face and Gesture Recognition (FG)*, pp. 1–7. IEEE, 2023.
- Nan Wu, Stanislaw Jastrzebski, Kyunghyun Cho, and Krzysztof J Geras. Characterizing and overcoming the greedy nature of learning in multi-modal deep neural networks. In *International Conference on Machine Learning*, pp. 24043–24055. PMLR, 2022.
- Yan Xia, Hai Huang, Jieming Zhu, and Zhou Zhao. Achieving cross modal generalization with multimodal unified representation. *Advances in Neural Information Processing Systems*, 36, 2024.
- Teng Xiao, Chao Cui, Huaisheng Zhu, and Vasant G. Honavar. Molbind: Multimodal alignment of language, molecules, and proteins, 2024. URL <https://arxiv.org/abs/2403.08167>.
- Ran Xu, Caiming Xiong, Wei Chen, and Jason Corso. Jointly modeling deep video and compositional text to bridge vision and language in a unified framework. In *Proceedings of the AAAI conference on artificial intelligence*, volume 29, 2015.
- Fengyu Yang, Chao Feng, Ziyang Chen, Hyounseob Park, Daniel Wang, Yiming Dou, Ziyao Zeng, Xien Chen, Rit Gangopadhyay, Andrew Owens, and Alex Wong. Binding touch to everything: Learning unified multimodal tactile representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 26340–26353, June 2024a.
- Fengyu Yang, Chao Feng, Daniel Wang, Tianye Wang, Ziyao Zeng, Zhiyang Xu, Hyounseob Park, Pengliang Ji, Hanbin Zhao, Yuanning Li, and Alex Wong. Neurobind: Towards unified multimodal representations for neural signals, 2024b. URL <https://arxiv.org/abs/2407.14020>.
- J Yang, Y Wang, R Yi, Y Zhu, A Rehman, A Zadeh, S Poria, and L-P Morency. Mtag: Modal-temporal attention graph for unaligned human multimodal language sequences. In *Proceedings of the 17th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL), 2021*, 2021.
- Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Mm-react: Prompting chatgpt for multimodal reasoning and action. *arXiv preprint arXiv:2303.11381*, 2023.

- Renrui Zhang, Ziyu Guo, Wei Zhang, Kunchang Li, Xupeng Miao, Bin Cui, Yu Qiao, Peng Gao, and Hongsheng Li. Pointclip: Point cloud understanding by clip. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8552–8562, 2022.
- Xiaohui Zhang, Jaehong Yoon, Mohit Bansal, and Huaxiu Yao. Multimodal representation learning by alternating unimodal adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 27456–27466, 2024.
- Jiazhou Zhou, Xu Zheng, Yuanhuiyi Lyu, and Lin Wang. E-clip: Towards label-efficient event-based open-world understanding by clip. *arXiv preprint arXiv:2308.03135*, 2023.
- Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiayi Cui, WANG HongFa, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, Cai Wan Zhang, Zhifeng Li, Wei Liu, and Li Yuan. Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=QmZKc7UZCy>.
- Lingyu Zhu and Esa Rahtu. V-slowfast network for efficient visual sound separation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1464–1474, 2022.

A RELATED WORK

A.1 MULTIMODAL LEARNING

Multimodal learning has gained significant attention in recent years due to its potential to enhance machine learning models by leveraging multiple data modalities, such as text, images, audio, and video. The fusion of these modalities aims to mimic human-like perception, improving performance in various applications ranging from healthcare to natural language processing. Popular supervised multimodal learning tasks include audio-visual classification (Peng et al., 2022; Feichtenhofer et al., 2019; Zhu & Rahtu, 2022), visual question answering (Antol et al., 2015; Guo et al., 2021), vision-language (Xu et al., 2015; Radford et al., 2021), vision-audio-language (Aytar et al., 2017; Harwath et al., 2018), and so on. In general, multimodal learning models fuse multiple unimodal features learned by unimodal encoders (Seichter et al., 2021; Nagrani et al., 2021; Wu et al., 2022; Wang et al., 2020a; Peng et al., 2022). For instance, Madaan et al. (2024) propose inter- and intra-modality modeling frameworks by considering the target as a source of multiple modalities. Du et al. (2023) propose to choose a targeted late-fusion learning method in supervised multi-modal tasks and prove insufficient learning of uni-modal features on each modality is negatively associated with the generalization-ability of the model. Zhang et al. (2024) decomposes the joint optimization into alternating unimodal learning scenarios by combining individual modality encoders with a shared head across all modalities.

A.2 MULTIMODAL ALIGNMENT

Multimodal learning addresses four key challenges (Liang et al., 2024c; Baltrušaitis et al., 2018; Liang et al., 2024d): managing interactions among redundant, unique, and synergistic features (Dumas et al., 2017; Liang et al., 2024a;b), aligning fine-grained and coarse-grained information (Wang et al., 2023; 2024a), reasoning across diverse features (Yang et al., 2023), and integrating external knowledge (Shen et al., 2022; Lyu et al., 2024). Among these challenges, multimodal alignment is one of the core challenges that many researchers aim to solve.

A common method in multimodal alignment is using cross-modal alignment by using attention mechanisms between pairwise modalities, such as vision-language (Tan & Bansal, 2019) and vision-language-audio (Tsai et al., 2019). Another effective approach is leveraging graph neural networks to align multimodal datasets (Yang et al., 2021; Wilf et al., 2023). For instance, Yang et al. (2021) transforms unaligned multimodal sequence data into nodes, with edges capturing interactions across modalities over time. Wilf et al. (2023) build graph structures for each modality—visual, textual, and acoustic—and create edges to represent their interactions.

To enhance the generalizability of cross-modal representations, Xia et al. (2024) employ a unified codebook approach, facilitating a joint embedding space for visual and audio modalities. Another prominent method (Radford et al., 2021) achieves cross-modal alignment by leveraging large collections of image-text pairs, making it a widely adopted strategy in multimodal learning (Zhang et al., 2022; Guzhov et al., 2022; Zhou et al., 2023).

A.3 BINDING METHODS

Recent studies have focused on aligning multimodal datasets by leveraging binding properties in various modalities. ImageBind (Girdhar et al., 2023) aligns multimodal data by using image representation as the anchor and aligning each modality embedding with the image embedding. Similarly, LanguageBind (Zhu et al., 2024) uses language representation as the anchor, aligning other modalities into the language space. PointBind (Guo et al., 2023) learns a joint embedding space across 3d point, language, image, and audio modalities by designating the point space as the central representation. Thanks to the efficacy of such a binding idea with a fixed anchor, several “-Bind” approaches have been studied in numerous domains (Teng et al., 2024; Xiao et al., 2024; Gao et al., 2024; Yang et al., 2024b; Balemans et al., 2024; Dhakal et al., 2024; Yang et al., 2024a). While these methods demonstrate strong performance in zero-shot cross-modality retrieval and classification tasks, they are constrained by their reliance on an existing single anchor modality.

Several approaches have integrated additional knowledge into multimodal representation spaces to address this limitation. Freebind (Wang et al., 2024a) introduces bi-modality spaces to enhance a

pretrained image-paired unified space. It generates pseudo-embedding pairs across diverse modality pairs and aligns them with the pre-trained unified space using contrastive learning. Omnibind (Wang et al., 2024b) leverages multiple pretrained multimodal models to construct pseudo item-pair retrievals based on top-1 recall across various modality combinations using pairwise cross-modal alignment. Both methods show promising results in cross-modal retrieval by incorporating extra spaces into existing pairwise binding spaces. However, they still rely on fixed (pre-trained) representation spaces.

Unibind (Lyu et al., 2024) highlights the imbalanced representation when using image-centered representation spaces. To address this, Unibind employs large language models (LLMs) to create a unified and balanced representation space. It constructs a knowledge base with multimodal category descriptions, establishes LLM-augmented class-wise embedding centers, and aligns other modalities to these centers through contrastive learning. This approach attempts to balance representations across modalities but still depends heavily on large-scale pretrained LLMs and centers alignment around a single unified space.

B PROOFS

B.1 PROOF OF PROPOSITION 1

Using the chain rule of the mutual information, we observe that

$$\begin{aligned} I(\mathbf{X}_1, f_1^{\text{suf}}(\mathbf{X}_1); \mathbf{X}_i) &= I(\mathbf{X}_1; \mathbf{X}_i) + I(f_1^{\text{suf}}(\mathbf{X}_1); \mathbf{X}_i | \mathbf{X}_1) \\ &= I(f_1^{\text{suf}}(\mathbf{X}_1); \mathbf{X}_i) + I(\mathbf{X}_1; \mathbf{X}_i | f_1^{\text{suf}}(\mathbf{X}_1)), \end{aligned} \quad (12)$$

Since $f_1^{\text{suf}}(\mathbf{X}_1)$ is a deterministic function of \mathbf{X}_1 , we have

$$I(f_1^{\text{suf}}(\mathbf{X}_1); \mathbf{X}_i | \mathbf{X}_1) = 0. \quad (13)$$

Moreover, f_1^{suf} obtained in Definition 1 with proper choice of \mathcal{Z} achieves the maximum mutual information, implying together with $I(\mathbf{X}; \mathbf{Y}) \leq \min\{H(\mathbf{X}), H(\mathbf{Y})\}$ that $I(f_1^{\text{suf}}(\mathbf{X}_1); \mathbf{X}_1) = H(\mathbf{X}_1)$, where $H(\mathbf{X}_1)$ is the entropy of \mathbf{X}_1 (Polyanskiy & Wu, 2024). In other words, we have $H(\mathbf{X}_1 | f_1^{\text{suf}}(\mathbf{X}_1)) = H(\mathbf{X}_1) - I(f_1^{\text{suf}}(\mathbf{X}_1); \mathbf{X}_1) = 0$. This gives

$$\begin{aligned} I(\mathbf{X}_1; \mathbf{X}_i | f_1^{\text{suf}}(\mathbf{X}_1)) &= H(\mathbf{X}_1 | f_1^{\text{suf}}(\mathbf{X}_1)) - H(\mathbf{X}_1 | f_1^{\text{suf}}(\mathbf{X}_1), \mathbf{X}_i) \\ &= 0 \end{aligned} \quad (14)$$

Substituting (13) and (14) into (12) yields

$$I(\mathbf{X}_1; \mathbf{X}_i) = I(f_1^{\text{suf}}(\mathbf{X}_1); \mathbf{X}_i). \quad (15)$$

We conclude the proof of Proposition 1 by noting that the optimality of FABIND (i.e., $I(f_1^{\text{suf}}(\mathbf{X}_1); \mathbf{X}_i) = I(f_1^{\text{suf}}(\mathbf{X}_1); f_i^{\text{FB}}(\mathbf{X}_i))$, $\forall i \in \{2, \dots, M\}$) yields

$$I(\mathbf{X}_1; \mathbf{X}_i) = I(f_1^{\text{suf}}(\mathbf{X}_1); f_i^{\text{FB}}(\mathbf{X}_i)). \quad (16)$$

B.2 PROOF OF PROPOSITION 2

Using the chain rule of mutual information, we have

$$\begin{aligned} I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_1, \mathbf{X}_i) &= I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_1) + I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_i | \mathbf{X}_1) \\ &= I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_i) + I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_1 | \mathbf{X}_i). \end{aligned} \quad (17)$$

Moreover, since $f_1^{\text{ins}}(\mathbf{X}_1)$ is a deterministic function of \mathbf{X}_1 , we have $I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_i | \mathbf{X}_1) = 0$, leading to $I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_1) = I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_i) + I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_1 | \mathbf{X}_i)$. Then, using the assumption $I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_1) < \epsilon$, it follows that

$$\begin{aligned} \epsilon &> I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_i) + I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_1 | \mathbf{X}_i) \\ &\stackrel{(a)}{\geq} I(f_1^{\text{ins}}(\mathbf{X}_1); \mathbf{X}_i) \\ &\stackrel{(b)}{\geq} I(f_1^{\text{ins}}(\mathbf{X}_1); f_i^{\text{FB}}(\mathbf{X}_i)), \end{aligned} \quad (18)$$

where the labeled inequalities follow from: (a) the non-negativity of mutual information; (b) the data processing inequality. This concludes the proof of Proposition 2.

B.3 PROOF OF THEOREM 1

To prove Theorem 1, we leverage the reverse inequality of M -variable Hölder inequality (Seo, 2013, eq. (2.8)). For the sake of completeness, we state the inequality in Lemma 1.

Lemma 1 (Reverse inequality of the M -variable Hölder inequality (Seo, 2013)). *Consider M sequences $(x_{i,j})_{j \in [n]}$, $i \in [M]$ of n positive scalars such that for some $0 < c_m \leq c_M < \infty$,*

$$0 < c_m \leq x_{i,j} \leq c_M < \infty, \forall i, j. \quad (19)$$

Then,

$$\prod_{i=1}^M \left(\sum_{j=1}^n x_{i,j} \right)^{\frac{1}{n}} \leq \frac{(c_m + c_M)^2}{4c_m c_M} \sum_{j=1}^n \left(\prod_{i=1}^M x_{i,j} \right)^{\frac{1}{n}}. \quad (20)$$

Now we start by writing the summation of InfoNCE losses for each $f_l^{(t)}(\mathbf{x}'_{l,k}), l \in [M]$ to $f_i(\mathbf{X}_i)$ as

$$\sum_{l=1}^M I_{\text{NCE}}(f_l(\mathbf{X}'_l); f_i(\mathbf{X}_i) | \tau) = -\frac{1}{|\mathcal{I}_B|} \sum_{k=1}^M \sum_{l=1}^M \log \frac{\exp\left(\frac{f_l^\top(\mathbf{x}'_{l,k}) f_i(\mathbf{x}_{i,k})}{\tau}\right)}{\sum_{j \in \mathcal{I}_B} \exp\left(\frac{f_l^\top(\mathbf{x}'_{l,k}) f_i(\mathbf{x}_{i,j})}{\tau}\right)}. \quad (21)$$

Then, the inner summation in (21) is bounded as

$$\begin{aligned} & \sum_{l=1}^M \log \frac{\exp\left(\frac{f_l^\top(\mathbf{x}'_{l,k}) f_i(\mathbf{x}_{i,k})}{\tau}\right)}{\sum_{j \in \mathcal{I}_B} \exp\left(\frac{f_l^\top(\mathbf{x}'_{l,k}) f_i(\mathbf{x}_{i,j})}{\tau}\right)} \\ &= \frac{1}{\tau} \sum_{l=1}^M f_l^\top(\mathbf{x}'_{l,k}) f_i(\mathbf{x}_{i,k}) - \log \prod_{l=1}^M \sum_{j \in \mathcal{I}_B} \exp\left(\frac{f_l^\top(\mathbf{x}'_{l,k}) f_i(\mathbf{x}_{i,j})}{\tau}\right) \\ &\stackrel{(a)}{\geq} \frac{1}{\tau} \sum_{l=1}^M f_l^\top(\mathbf{x}'_{l,k}) f_i(\mathbf{x}_{i,k}) - \log \left(C_{\mathcal{F},k,i} \sum_{j \in \mathcal{I}_B} \prod_{l=1}^M \exp\left(\frac{f_l^\top(\mathbf{x}'_{l,k}) f_i(\mathbf{x}_{i,j})}{\tau |\mathcal{I}_B|}\right) \right)^{|\mathcal{I}_B|} \\ &\stackrel{(b)}{=} \frac{M}{\tau} \mathbf{a}_k^\top f_i(\mathbf{x}_{i,k}) - |\mathcal{I}_B| \log \sum_{j \in \mathcal{I}_B} \exp\left(\frac{M \mathbf{a}_k^\top f_i(\mathbf{x}_{i,j})}{\tau |\mathcal{I}_B|}\right) - |\mathcal{I}_B| \log C_{\mathcal{F},k,i} \\ &= |\mathcal{I}_B| \log \exp\left(\frac{M \mathbf{a}_k^\top f_i(\mathbf{x}_{i,k})}{\tau |\mathcal{I}_B|}\right) - |\mathcal{I}_B| \log \sum_{j \in \mathcal{I}_B} \exp\left(\frac{M \mathbf{a}_k^\top f_i(\mathbf{x}_{i,j})}{\tau |\mathcal{I}_B|}\right) - |\mathcal{I}_B| \log C_{\mathcal{F},k,i} \\ &= |\mathcal{I}_B| \log \frac{\exp\left(\frac{M \mathbf{a}_k^\top f_i(\mathbf{x}_{i,k})}{\tau |\mathcal{I}_B|}\right)}{\sum_{j \in \mathcal{I}_B} \exp\left(\frac{M \mathbf{a}_k^\top f_i(\mathbf{x}_{i,j})}{\tau |\mathcal{I}_B|}\right)} - |\mathcal{I}_B| \log C_{\mathcal{F},k,i}, \end{aligned} \quad (22)$$

where the labeled (in)equalities follow from: (a) Lemma 1 and $C_{\mathcal{F},k,i} = \frac{(c_{\mathcal{F},k,i}^{\min} + c_{\mathcal{F},k,i}^{\max})^2}{4c_{\mathcal{F},k,i}^{\min} c_{\mathcal{F},k,i}^{\max}}$ with

$$\begin{aligned} c_{\mathcal{F},k,i}^{\min} &= \min_{\ell \in [M], j \in \mathcal{I}_B} \exp\left(\frac{f_\ell^\top(\mathbf{x}'_{l,k}) f_i(\mathbf{x}_{i,j})}{\tau}\right), \text{ and} \\ c_{\mathcal{F},k,i}^{\max} &= \max_{\ell \in [M], j \in \mathcal{I}_B} \exp\left(\frac{f_\ell^\top(\mathbf{x}'_{l,k}) f_i(\mathbf{x}_{i,j})}{\tau}\right); \end{aligned} \quad (23)$$

and (b) the definition of anchor embedding (7). Substituting (22) into (21) gives

$$\begin{aligned} \sum_{l=1}^M I_{\text{NCE}}(f_l(\mathbf{X}'_l); f_i(\mathbf{X}_i) | \tau) &\leq -\frac{1}{|\mathcal{I}_B|} \sum_{k=1}^M \left[|\mathcal{I}_B| \log \frac{\exp\left(\frac{M \mathbf{a}_k^\top f_i(\mathbf{x}_{i,k})}{\tau |\mathcal{I}_B|}\right)}{\sum_{j \in \mathcal{I}_B} \exp\left(\frac{M \mathbf{a}_k^\top f_i(\mathbf{x}_{i,j})}{\tau |\mathcal{I}_B|}\right)} - |\mathcal{I}_B| \log C_{\mathcal{F},k,i} \right] \\ &= |\mathcal{I}_B| I_{\text{NCE}}\left(\mathbf{A}; f_i(\mathbf{X}_i) \middle| \frac{\tau |\mathcal{I}_B|}{M}\right) + \sum_{k=1}^M \log C_{\mathcal{F},k,i}. \end{aligned} \quad (24)$$

Rearranging (24) and setting $\tilde{\tau} = \frac{\tau|\mathcal{I}_B|}{M}$ in (23) and (24) yield

$$I_{\text{NCE}}(\mathbf{A}; f_i(\mathbf{X}_i) \mid \tilde{\tau}) \geq \frac{1}{|\mathcal{I}_B|} \sum_{l=1}^M I_{\text{NCE}} \left(f_l(\mathbf{X}'_l); f_i(\mathbf{X}_i) \mid \frac{\tilde{\tau}M}{|\mathcal{I}_B|} \right) - \frac{1}{|\mathcal{I}_B|} \sum_{k=1}^{|\mathcal{I}_B|} \log C_{\mathcal{F},k,i}, \quad (25)$$

which concludes the proof of Theorem 1.

C EXPERIMENT DETAILS

C.1 EXPERIMENTS WITH SYNTHETIC DATASETS

Synthetic datasets. We employ a latent variable model (Bishop & Nasrabadi, 2006) for generating synthetic multimodal datasets. A latent variable model is a statistical model for data $\mathbf{X} \in \mathbb{R}^{d_x}$, under which \mathbf{X} is generated according to a conditional probability distribution $P_{\mathbf{X}|\mathbf{Z}}$, where $\mathbf{Z} \in \mathbb{R}^{d_z}$ is the latent variable. In terms of the representation learning framework, \mathbf{Z} can be seen as a true representation of \mathbf{X} . Moreover, we assume that the class label $\mathbf{Y} \in [K]$ and the latent variable \mathbf{Z} are jointly distributed according to $P_{\mathbf{Z},\mathbf{Y}}$.

For the marginal distribution of \mathbf{Z} , we make use of a Gaussian mixture model (GMM) (Bishop & Nasrabadi, 2006), and hence the probability density function (PDF) of \mathbf{Z} is a weighted sum of Gaussian densities. In particular, the PDF of \mathbf{Z} is defined as follows:

$$p_{\mathbf{Z}}(\mathbf{z}) = \prod_{y=1}^K \pi_y \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y), \quad (26)$$

where K is the number of mixture components, $\pi_y = \Pr(\mathbf{Y} = y)$ is the component prior probability, and $\mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$ denotes Gaussian PDF with mean $\boldsymbol{\mu}_y \in \mathbb{R}^{d_z}$ and covariance matrix $\boldsymbol{\Sigma}_y \in \mathbb{R}^{d_z \times d_z}$. This leads to the conditional PDF of \mathbf{Z} as $p_{\mathbf{Z}|\mathbf{Y}}(\mathbf{z}|y) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$.

Once a latent variable \mathbf{z} is generated from GMM in (26), we generate data samples $(\mathbf{x}_{i,1}, \mathbf{x}_{i,2}, \dots, \mathbf{x}_{i,N})$ for i -th modality using the conditional PDFs of \mathbf{X}_i given \mathbf{z} , denoted by $p_{\mathbf{X}_i|\mathbf{Z}}(\mathbf{x}_i|\mathbf{z})$. Specifically, we use the model $\mathbf{X}_i = g_i(\mathbf{Z}_i) + \mathbf{N}$, where $g_i : \mathbb{R}^{d_z} \rightarrow \mathbb{R}^{d_x}$ is a non-linear projection from latent space to observation space, and $\mathbf{N} \sim \mathcal{N}(\mathbf{0}, I_{d_x})$ is Gaussian noise with zero-mean and identity covariance matrix. To make the inherent correlation between \mathbf{X}_i and \mathbf{Z}_i different among modalities, we choose g_i such that

$$g_i(\mathbf{Z}) = \Theta_i^{(2)} \text{sigmoid} \left(\Theta_i^{(1)} \mathbf{Z} \right), \quad (27)$$

where $\text{sigmoid}(x) = \frac{1}{1+e^{-x}}$ is applied element-wise, and $\Theta_i^{(1)} \in \mathbb{R}^{d_x \times d_z}$ and $\Theta_i^{(2)} \in \mathbb{R}^{d_x \times d_x}$ are matrices randomly generated from Gaussian distribution. Moreover, after $\Theta_i^{(1)}, i \in [M]$ are generated, we set arbitrary columns of them all zero, so that the number of all zero columns decreases in i . For example, 60% of columns of $\Theta_1^{(1)}$ are all-zero, while only 10% of columns of $\Theta_M^{(1)}$ are all-zero. This enables approximate control the correlation between \mathbf{X}_i and \mathbf{Z} , providing estimates of best modality (\mathbf{X}_M) or worst modality (\mathbf{X}_1). To have meaningful labels for this latent model, which requires for downstream tasks, we set the labels \mathbf{Y} being the component index in GMM. In particular, since there are K components in GMM (26), there exists K categories in \mathbf{Y} . We conduct experiments with three different synthetic datasets by setting $M = 4, 6, 8$. For all synthetic datasets, we fix $d_x = 16$, $d_z = 8$, and $K = 50$.

Experiment details. We initialize two different versions of backbones for all modalities, where the first is a random backbone (highlighted by (rnd) in figures), and the second is a backbone pretrained with InfoNCE loss. For each backbone, we use a simple multilayer perceptron (MLP). Comparing the results with these two versions of backbone provides how much both FABIND and CENTROBIND are robust to backbone quality. Given the backbones for M modalities, we align the corresponding embedding spaces using either FABIND with anchor \mathbf{X}_i (denoted by \mathbf{X}_i -B in figures) or CENTROBIND (denoted by CB in figures). Finally, with the encoders aligned by either FABIND or CENTROBIND, we evaluate classification accuracy as a measure of representation quality. We use a simple MLP for the classifier. To distinguish between accuracy with embeddings from a single

modality and the one with concatenated embeddings from all modalities, we denote by $\text{acc}(\mathbf{Z}_i)$ the accuracy with embeddings from i -th modality and by $\text{acc}(\text{All})$ the accuracy with embeddings from all modalities.

Representation visualization. Figure 5 shows t-SNE (Van der Maaten & Hinton, 2008) and UMAP (McInnes et al., 2018) visualization of embeddings from FABIND and CENTROBIND. In both t-SNE and UMAP visualizations, CENTROBIND yields more clustered representations, while FABIND constructs scattered representation. This validates our analysis that CENTROBIND creates better representation space due to the ability to learn intra- and shared information.

Additional experiments with different number of modalities. Additional experimental results on synthetic datasets with $M = 6, 8$ number of modalities are shown in Figure 6 and Figure 7. These experiments verify that CENTROBIND is able to handle large number of modalities.

Comparison with other dynamic anchor generation. We compare the centroid dynamic anchor method with other possible approaches such as weighted average and random anchor fixing. In Figure 8, we plot the accuracies of each methods for special cases, where modalities are unevenly distributed. Specifically, we create 4 modalities with each modal quality strictly different. For the experiments (a) and (b) in Figure 8, we set \mathbf{X}_1 , \mathbf{X}_2 , and \mathbf{X}_3 to be very uninformative, while \mathbf{X}_4 being high-quality dataset. The experiments for (c) and (d) in Figure 8 is conducted with poor quality datasets of \mathbf{X}_1 and \mathbf{X}_2 and high-quality dataset of \mathbf{X}_3 and \mathbf{X}_4 . For the weighted average (denoted as WAB in Figure 8), we choose weight depending on the quality, (0.2, 0.2, 0.2, 1) for (a) and (b) experiments and (0.2, 0.2, 0.8, 0.8) for (c) and (d) experiments, which are identical to the information rate of each modalities. For the random modality for dynamic anchor (denoted as RB in Figure 8), we randomly choose one of modalities for dynamic anchor at each iteration. Note that we freeze the anchor encoder for the random anchor. To see how impact the intra modal learning, we conduct random anchor with intra-information learning (denoted as RB+Intra), for which we randomly choose anchor modality at each iteration and we skip the freezing the encoder so that all encoders (including the anchor encoder) are trained.

Such scenario where we differentiate the modal distribution is in general called *modality imbalance* problem (Du et al., 2023; Peng et al., 2022; Zhang et al., 2024). Intuitively, if modality imbalance problem exists, centroid may lead to suboptimal construction of dynamic anchor, and other methods (e.g., weighted average) performs better. Nevertheless, CENTROBIND performs better or similar to the weighted average methods, showing its robustness to the modality imbalance problem. Moreover, we observe that the random anchor method also performs moderately. Through the experiments, we conjecture that dynamic anchor generation method is not much important for the final performance whenever every encoder is trained during training.

The modality imbalance problem requires other information, such as domain knowledge, labels, or downstream tasks. Since this work mainly focuses on multimodal alignment under contrastive learning, for which we do not assume such information is available, and hence we leave the modality imbalance problem for dynamic anchor as a future work.

Convergence and stability analysis. Convergence rate of CENTROBIND might be different than the one of FABIND as we replaced the fixed anchor with dynamic anchor. We plot the loss curves of CENTROBIND and FABIND during training in Figure 4. It shows that the loss of CENTROBIND is saturated earlier than that of FABIND. We think that this is because centroid is a minimizer of embeddings in terms of Euclidean distance, and hence converging the embeddings to their centroid is easier than converging to one of embeddings. The plot further shows that there exists a cross point where loss curves are across each other. We think that this happens due to the number of InfoNCE losses optimized by CENTROBIND and FABIND. Specifically, with M modalities, CENTROBIND minimizes M InfoNCE losses, whereas FABIND minimizes $M - 1$ InfoNCE losses. This leads to smaller loss from FABIND when the encoders are well-trained, indicating the cross point in Figure 4.

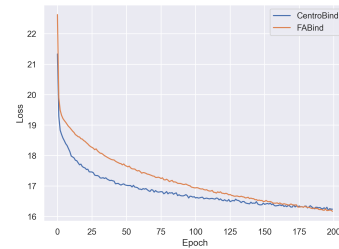


Figure 4: Training loss.

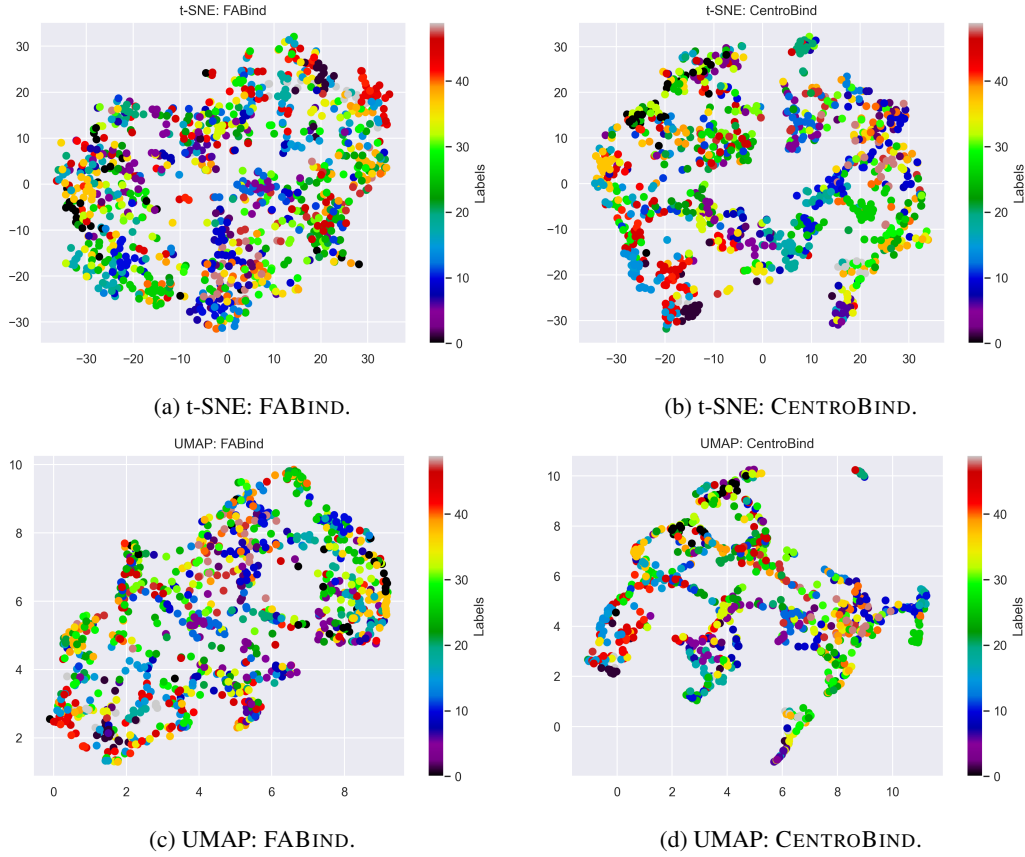


Figure 5: Representation visualization via t-SNE and UMAP.

C.2 EXPERIMENTS WITH REAL-WORLD DATASETS

Training details. We utilize Low-Rank Adaptation (Hu et al., 2022) for training CENTROBIND and FABIND, enhancing training efficiency and achieving impressive results with fewer iterations. For parameter settings, we set a learning rate of 0.001, the AdamW optimizer (Loshchilov & Hutter, 2019) with a batch size of 16, and a temperature of 0.3 for InfoNCE. Training CENTROBIND requires augmentation. We augment video frames with various transformations, including random perspective shifts, random flips and rotation, color jitter, Gaussian blur, and auto-contrast adjustment. For the audio modality, we apply a low-pass filter, speed changes, echo effect, room impulse response convolution, and background noise. For the text modality, we generate paraphrased sentences using the Phi-3 language model served using Ollama⁴.

⁴<https://ollama.com/library/phi3>

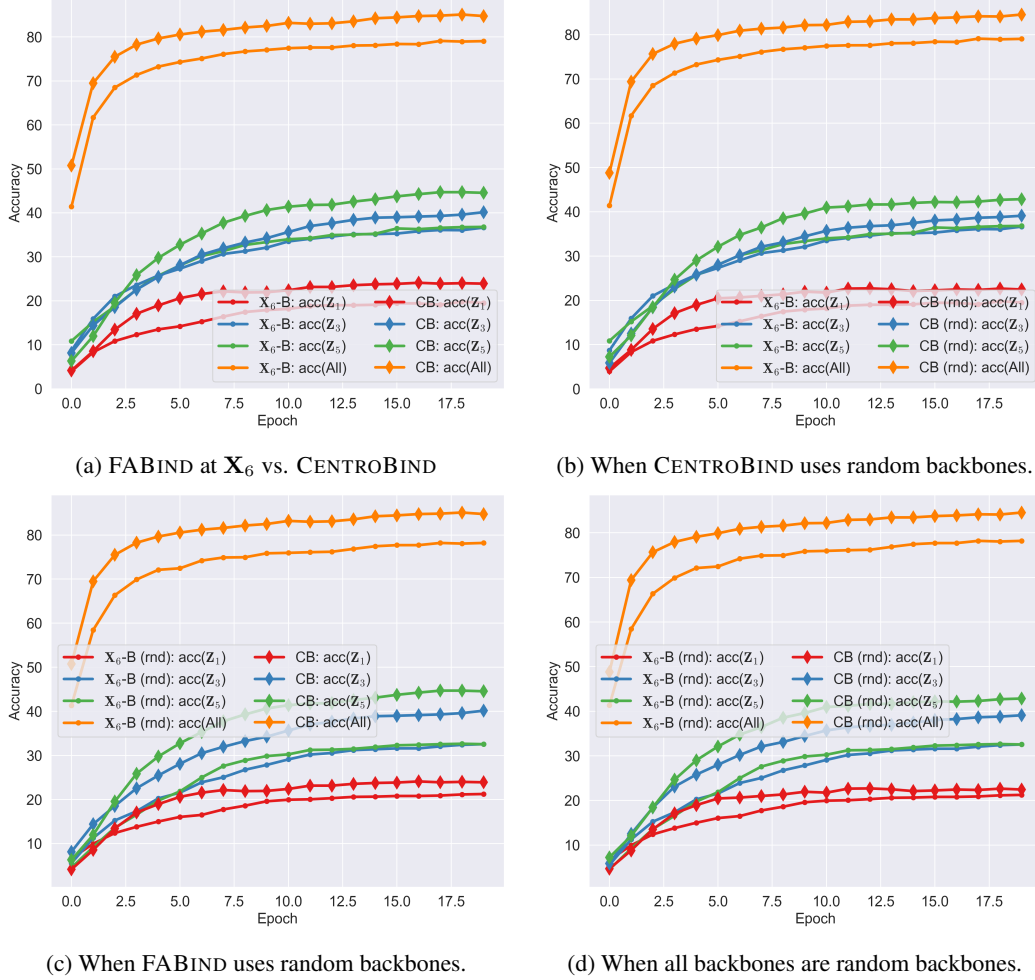


Figure 6: Experiment results with synthetic dataset of $M = 6$ modalities. Abbreviation: X_i -B or CB: applying FABIND method to backbones with anchor X_i or applying CENTROBIND; acc(Z_i) or acc(All): accuracy of Z_i or of concatenated embeddings (Z_1, \dots, Z_M); (rnd): if random backbones are used for X_i -B or CB.

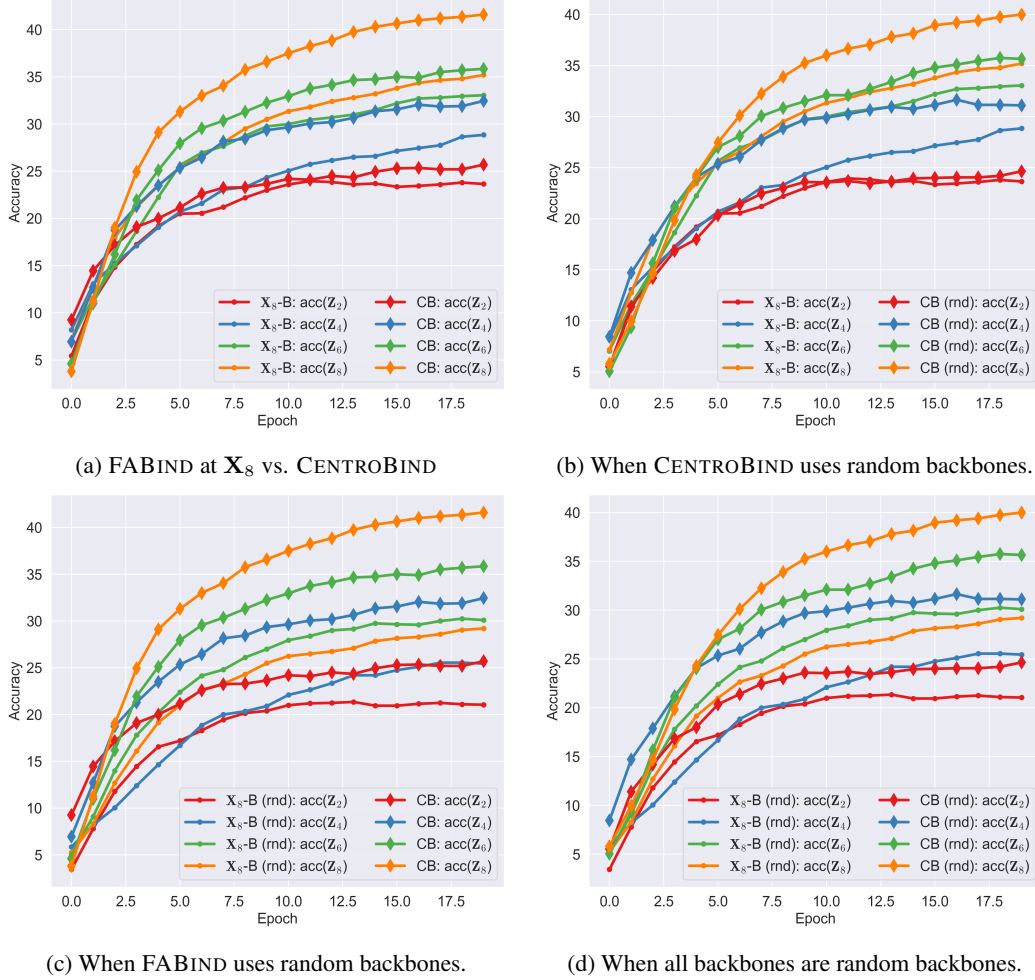


Figure 7: Experiment results with synthetic dataset of $M = 8$ modalities. Abbreviation: \mathbf{X}_i -B or CB: applying FABIND method to backbones with anchor \mathbf{X}_i or applying CENTROBIND; $\text{acc}(\mathbf{Z}_i)$ or $\text{acc}(\text{All})$: accuracy of \mathbf{Z}_i or of concatenated embeddings $(\mathbf{Z}_1, \dots, \mathbf{Z}_M)$; (rnd): if random backbones are used for \mathbf{X}_i -B or CB.

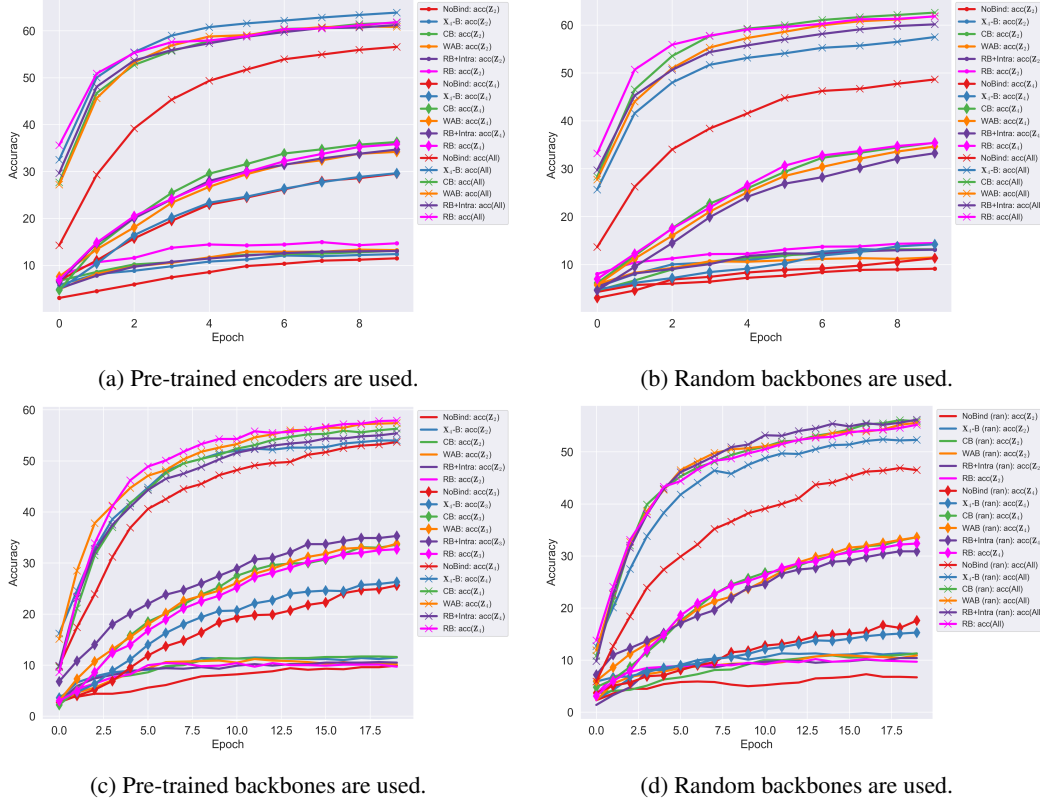


Figure 8: Comparison of other dynamic anchor generation methods. (a) and (b): Modal qualities are set to (0.2, 0.2, 0.2, 1). (c) and (d): Modal qualities are set to (0.2, 0.2, 0.8, 0.8). Abbreviation: X_i -B or CB: applying FABIND method to backbones with anchor X_i or applying CENTROBIND; WAB: weighted average for dynamic anchor with weight identical to the predefined quality for each modality; RB+Intra: randomly choosing a modality for dynamic anchor in every iteration and intra information learning; RB: randomly choosing a modality for dynamic anchor in every iteration; $\text{acc}(\mathbf{Z}_i)$ or $\text{acc}(\text{All})$: accuracy of \mathbf{Z}_i or of concatenated embeddings $(\mathbf{Z}_1, \dots, \mathbf{Z}_M)$; (ran): if random backbones are used.