

Contrastive Multimodal Fusion with TupleInfoNCE

Yunze Liu^{1,7} Qingnan Fan³ Shanghang Zhang⁴ Hao Dong^{5,6,8} Thomas Funkhouser² Li Yi^{1,2*}

¹IIS, Tsinghua University ²Google Research ³Stanford University
⁴UC Berkeley ⁵CFCS, CS Dept., Peking University ⁶AiIT, Peking University
⁷Xidian University ⁸Peng Cheng Laboratory

{liuyczchina, fqnchina}@gmail.com, {eric yi, tfunkhouser}@google.com
shz@eecs.berkeley.edu, hao.dong@pku.edu.cn

Abstract

This paper proposes a method for representation learning of multimodal data using contrastive losses. A traditional approach is to contrast different modalities to learn the information shared among them. However, that approach could fail to learn the complementary synergies between modalities that might be useful for downstream tasks. Another approach is to concatenate all the modalities into a tuple and then contrast positive and negative tuple correspondences. However, that approach could consider only the stronger modalities while ignoring the weaker ones. To address these issues, we propose a novel contrastive learning objective, TupleInfoNCE. It contrasts tuples based not only on positive and negative correspondences, but also by composing new negative tuples using modalities describing different scenes. Training with these additional negatives encourages the learning model to examine the correspondences among modalities in the same tuple, ensuring that weak modalities are not ignored. We provide a theoretical justification based on mutual-information for why this approach works, and we propose a sample optimization algorithm to generate positive and negative samples to maximize training efficacy. We find that TupleInfoNCE significantly outperforms previous state of the arts on three different downstream tasks.

1. Introduction

Human perception of the world is naturally multimodal. What we see, hear, and feel all contain different kinds of information. Various modalities complement and disambiguate each other, forming a representation of the world. Our goal is to train machines to fuse such multimodal inputs to produce such representations in a self-supervised manner without manual annotations.

*Corresponding author

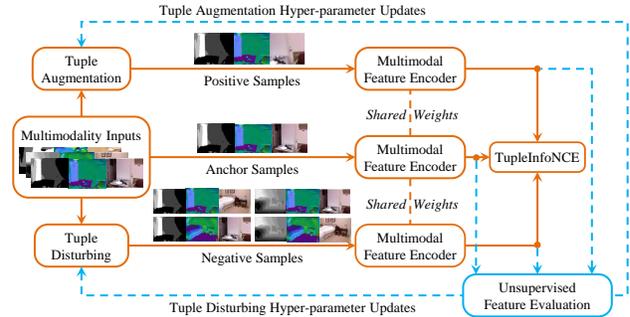


Figure 1. Overview of sample-optimized TupleInfoNCE.

An increasingly popular self-supervised representation learning paradigm is contrastive learning, which learns feature representations via optimizing a contrastive loss and solving an instance discrimination task [24, 12, 5]. Recently several works have explored contrastive learning for multimodal representation learning [31, 1, 21]. Among them, the majority [31, 1] learn a crossmodal embedding space – they contrast different modalities to capture the information shared across modalities. However, they do not examine the fused representation of multiple modalities directly, failing to fully leverage multimodal synergies. To cope with this issue, [21] proposes an RGB-D representation learning framework to directly contrast pairs of point-pixel pairs. However, it is restricted to two modalities only.

Instead of contrasting different data modalities, we propose to contrast multimodal input tuples, where each tuple element corresponds to one modality. We learn representations so that tuples describing the same scene (set of multimodal observations) are brought together while tuples from different scenes are pushed apart. This is more general than crossmodal contrastive learning. It not only supports extracting the shared information across modalities, but also allows modalities to disambiguate each other and to keep their specific information, producing better-fused representations.

However, contrasting tuples is not as straightforward as

contrasting single elements, especially if we want the learned representation to encode the information from each element in the tuple and to fully explore the synergies among them. The core challenge is: “which tuple samples to contrast?” Previously researchers [37, 21] have observed that always contrasting tuples containing corresponding elements from the same scene can converge to a lazy suboptimum where the network relies only on the strongest modality for scene discrimination. Therefore to avoid weak modalities being ignored and to facilitate modality fusion, we need to contrast from more challenging negative samples. Moreover, we need to optimize the positive samples as well so that the contrastive learning can keep the shared information between positive and anchor samples while abstracting away nuisance factors. Strong variations between the positive and anchor samples usually result in smaller shared information but a greater degree of invariance against nuisance variables. Thus a proper tradeoff is needed.

To handle the above challenges, we propose a novel contrastive learning objective named TupleInfoNCE (Figure 1). Unlike the popular InfoNCE loss [24], TupleInfoNCE is designed explicitly to facilitate multimodal fusion. TupleInfoNCE leverages positive samples generated via augmenting anchors and it exploits challenging negative samples whose elements are not necessarily in correspondence. These negative samples encourage a learning model to examine the correspondences among elements in an input tuple, ensuring that weak modalities and the modality synergy are not ignored. To generate such negative samples we present a tuple disturbing strategy with a theoretical basis for why it helps.

TupleInfoNCE also introduces optimizable hyper-parameters to control both the negative sample and the positive sample distributions. This allows optimizing samples through a hyper-parameter optimization process. We define reward functions regarding these hyper-parameters and measure the quality of learned representations via unsupervised feature evaluation. We put unsupervised feature evaluation in an optimization loop that updates these hyper-parameters to find a sample-optimized TupleInfoNCE (Figure 1).

We evaluate TupleInfoNCE on a wide range of multimodal fusion tasks including multimodal semantic segmentation on NYUv2 [29], multimodal object detection on SUN RGB-D [30] and multimodal sentiment analysis on CMU-MOSI [38] and CMU-MOSEI [39]. We demonstrate significant improvements over previous state-of-the-art multimodal self-supervised representation learning methods (+4.7 mIoU on NYUv2, +1.2 mAP@0.25 on SUN RGB-D, +1.0% acc7 on MOSI, and +0.5% acc7 on MOSEI).

Our key contributions are threefold. First, we present a novel TupleInfoNCE objective for contrastive multimodal fusion with a theoretical justification. Secondly, we pose the problem of optimizing TupleInfoNCE with a self-supervised approach to select the contrastive samples. Finally, we

demonstrate state-of-the-art performance on a wide range of multimodal fusion benchmarks and provide ablations to evaluate the key design decisions.

2. Related Work

2.1. Self-Supervised Multimodal Learning

Self-supervised learning (SSL) uses auxiliary tasks to learn data representation from the raw data without using additional labels [33, 14, 23, 10, 34], helping to improve the performance of the downstream tasks. Recently, research on SSL leverages multimodal properties of the data [8, 2, 31, 11, 1, 21]. The common strategy is to explore the natural correspondences among different views and use contrastive learning (CL) to learn representations by pushing views describing the same scene closer, while pushing views of different scenes apart [8, 2, 31, 11, 1]. We refer to this line of methods as crossmodal embedding, which focuses on extracting the information shared across modalities rather than examining the fused representation directly, failing to fully explore the modality synergy for multimodal fusion.

2.2. Contrastive Representation Learning

CL is a type of SSL that has received increasing attention for it brings tremendous improvements on representation learning. According to the learning method, it can be grouped into Instance-based [24, 12, 6, 5] and Prototype-based CL [17, 4]; According to the modality of data, it can be categorized into single-modality based [7, 13] and multi-modality based CL [1, 21, 31]. An underexplored challenge for CL is how to select hard negative samples to build the negative pair [15, 27, 13, 7]. Most existing methods either increase batch size or keep large memory banks, leading to large memory requirements [12]. Recently, several works study CL from the perspective of mutual information (MI). [32] argues MI between views should be reduced by data augmentation while keeping task-relevant information intact. [36] shows the family of CL algorithms maximizes a lower bound on MI between multi-“views” where typical views come from image augmentations, and finds the choice of negative samples and views are critical to these algorithms. We build upon this observation with an optimization framework for selecting contrastive samples.

2.3. AutoML

AutoML is proposed to automatically create models that outperform the manual design. The progress of neural architectural search (NAS) [40, 20, 3], data augmentation strategy search [9, 18] and loss function search [16] have greatly improved the performance of neural networks. But most of these methods focus on a supervised learning setting. Recently, developing AutoML techniques in an unsupervised/self-supervised learning scenario has drawn

more attention [19, 32, 22]. UnNAS [19] shows the potential of searching for better neural architectures with self-supervision. InfoMin [32] and SelfAugment [22] explore how to search better data augmentation for CL on 2D images. In our work, we focus on optimizing two key components of a multimodal CL framework unsupervisedly - data augmentation and negative sampling strategies, none of which has been previously explored for generic multimodal inputs.

3. Revisiting InfoNCE

Before describing our method, we first review the InfoNCE loss widely adopted for contrastive representation learning [24], and then discuss its limitations for multimodal inputs. Given an anchor random variable $\mathbf{x}_{1,i} \sim p(\mathbf{x}_1)$, the popular contrastive learning framework aims to differentiate a positive sample $\mathbf{x}_{2,i} \sim p(\mathbf{x}_2|\mathbf{x}_{1,i})$ from negative samples $\mathbf{x}_{2,j} \sim p(\mathbf{x}_2)$. This is usually done by minimizing the InfoNCE loss:

$$\mathcal{L}_{\text{NCE}} = -\mathbb{E} \left[\log \frac{f(\mathbf{x}_{2,i}, \mathbf{x}_{1,i})}{\sum_{j=1}^N f(\mathbf{x}_{2,j}, \mathbf{x}_{1,i})} \right] \quad (1)$$

where $f(\mathbf{x}_{2,j}, \mathbf{x}_{1,i})$ is a positive scoring function usually chosen as a log-bilinear model. It has been shown that minimizing \mathcal{L}_{NCE} is equivalent to maximizing a lower bound of the mutual information $I(\mathbf{x}_2; \mathbf{x}_1)$. Many negative samples are required to properly approximate the negative distribution $p(\mathbf{x}_2)$ and tighten the lower bound.

In the problem setting of multimodal inputs, an input sample can be represented as a K -tuple $\mathbf{t} = (\mathbf{v}^1, \mathbf{v}^2, \dots, \mathbf{v}^K)$ where each element \mathbf{v}^k corresponds to one modality and K denotes the total number of modalities being considered. A straightforward way of learning multimodal representations is to draw anchor samples $\mathbf{t}_{1,i} \sim p(\mathbf{t}_1)$, their positive samples $\mathbf{t}_{2,i} \sim p(\mathbf{t}_2|\mathbf{t}_{1,i})$ and negative samples $\mathbf{t}_{2,j} \sim p(\mathbf{t}_2)$, and then optimize the InfoNCE objective. However, previous works [37, 21] observe that even when $K = 2$ simply drawing negative samples from the marginal distribution $p(\mathbf{t}_2)$ is insufficient for learning good representations. Weak modalities tend to be largely ignored and synergies among modalities are not fully exploited. The issue becomes more severe when $K > 2$ when the informativeness of different modalities varies a lot.

Figure 2 provides an intuitive explanation. When one modality \mathbf{v}^k is particularly informative compared with the rest modalities $\bar{\mathbf{v}}^k$ in the input tuple \mathbf{t} , namely $I(\mathbf{v}_2^k; \mathbf{v}_1^k) \gg I(\bar{\mathbf{v}}_2^k; \bar{\mathbf{v}}_1^k)$, maximizing a lower bound of $I(\mathbf{t}_2; \mathbf{t}_1) = I(\mathbf{v}_2^k; \bar{\mathbf{v}}_2^k; \mathbf{v}_1^k; \bar{\mathbf{v}}_1^k)$ will be largely dominated by the modality specific information $I(\mathbf{v}_2^k; \mathbf{v}_1^k | \bar{\mathbf{v}}_2^k, \bar{\mathbf{v}}_1^k)$, which is usually not as important as the information shared across modalities $I(\mathbf{v}_2^k; \bar{\mathbf{v}}_2^k; \mathbf{v}_1^k; \bar{\mathbf{v}}_1^k)$. Overemphasizing the modality specific information from the strong modality might sacrifice the weak modalities and the modality synergy during learning.

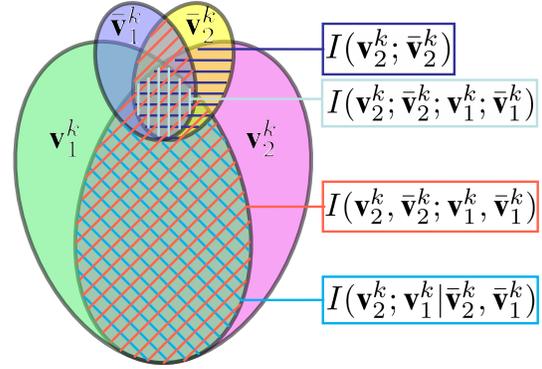


Figure 2. Information diagram

4. TupleInfoNCE

To alleviate the limitations of InfoNCE for overlooking weak modalities and the modality synergy, we present a novel TupleInfoNCE objective. We leverage a *tuple disturbing* strategy to generate challenging negative samples, which prevents the network from being lazy and only focusing on strong modalities. In addition, we introduce optimizable data augmentations which are applied to anchor samples for positive sample generation. We optimize both the positive and negative samples to balance the information contributed by each modality. All these are incorporated into the proposed TupleInfoNCE objective, designed explicitly to facilitate multimodal fusion.

4.1. Tuple disturbing and augmentation

Tuple disturbing Generating challenging negative samples is fundamentally important to learning effective representation in contrastive learning, especially in the case of multimodal fusion setting where the strong modalities tend to dominate the learned representation [21, 37]. We present a *tuple disturbing* strategy to generate negative samples where not all modalities are in correspondence and certain modalities exhibit different scenes.

Given an anchor sample $(\mathbf{v}_{1,i}^1, \dots, \mathbf{v}_{1,i}^k, \dots, \mathbf{v}_{1,i}^K)$ and its positive sample $(\mathbf{v}_{2,i}^1, \dots, \mathbf{v}_{2,i}^k, \dots, \mathbf{v}_{2,i}^K)$, we propose a k -disturbed negative sample represented as $(\mathbf{v}_{2,j}^1, \dots, \mathbf{v}_{2,d(j)}^k, \dots, \mathbf{v}_{2,j}^K)$, where $d(\cdot)$ is a disturbing function producing a random index from the sample set. The negative sample has $K - 1$ modalities $\bar{\mathbf{v}}_{2,j}^k$ from one scene and one modality $\mathbf{v}_{2,d(j)}^k$ from a different scene. Therefore, in order to correctly discriminate the positive sample from k -disturbed negative samples, the learned representation has to encode the information of the k -th modality, since the K -tuple could become negative only due to differences in the k -th modality. k -disturbed negative samples become especially challenging when they are only partially negative, e.g. $\bar{\mathbf{v}}_{2,j}^k$ becomes very similar to $\bar{\mathbf{v}}_{2,i}^k$. Simply treating \mathbf{v}^k as an independent modality without considering its correlation with the rest modalities is not able to fully suppress the score

of such partially negative samples in a log-bilinear model. Only when the network tells the disturbed modality $\mathbf{v}_{2,d(j)}^k$ is not in correspondence with the rest modalities $\bar{\mathbf{v}}_{2,j}^k$, can it fully suppress the partially negative samples. Therefore k -disturbed negative samples encourage the correlation between each modality and the rest to be explored.

We disturb each modality separately and generate K types of negative samples to augment the vanilla InfoNCE objective. This enforces the representation learning of each specific modality in the multimodal inputs. We use α_k to represent the ratio of k -disturbed negative samples. Intuitively, the larger α_k we use, the more emphasis we put on the k -th modality.

Tuple augmentation Given an anchor sample \mathbf{t}_1 , we apply the data augmentation to each modality separately to generate the positive sample \mathbf{t}_2 . The data augmentation applied to modality \mathbf{v}^k will directly influence $I(\mathbf{v}_2^k; \mathbf{v}_1^k)$ [32], which roughly measures the information contribution of modality \mathbf{v}^k in $I(\mathbf{t}_2; \mathbf{t}_1)$. To further balance the contribution of each modality in our fused representation, we parameterize these data augmentations with a hyper-parameter β and make β optimizable for different modalities.

4.2. Objective function

The TupleInfoNCE objective is designed for fusing the multimodal input tuple $\mathbf{t} = (\mathbf{v}^1, \mathbf{v}^2, \dots, \mathbf{v}^K)$. Given an anchor sample $\mathbf{t}_{1,i} \sim p(\mathbf{t}_1)$, we draw its positive sample $\mathbf{t}_{2,i} \sim p_\beta(\mathbf{t}_2|\mathbf{t}_{1,i})$, and negative sample $\mathbf{t}_{2,j|j \neq i} \sim q_\alpha(\mathbf{t}_2)$ following a ‘‘proposal’’ distribution where either all modalities are in correspondence yet stem from a different scene, or each modality is disturbed to encourage modality synergy. To be specific, with probability α_0 we sample negative samples from $p(\mathbf{t}_2)$, and with probability α_k we sample k -disturbed negative samples from $p(\bar{\mathbf{v}}_2^k)p(\mathbf{v}_2^k)$, where $\{\alpha_k\}_{k=0}^K$ is a set of prior probabilities balancing different types of negative samples which sum to 1. This essentially changes our negative sample distribution to be $q_\alpha(\mathbf{t}_2) = \alpha_0 p(\mathbf{t}_2) + \sum_{k=1}^K \alpha_k p(\bar{\mathbf{v}}_2^k)p(\mathbf{v}_2^k)$. Therefore, the TupleInfoNCE objective is defined as below:

$$\mathcal{L}_{\text{TNCE}}^{\alpha\beta} = - \mathbb{E}_{\substack{\mathbf{t}_{2,i} \sim p_\beta(\mathbf{t}_2|\mathbf{t}_{1,i}) \\ \mathbf{t}_{2,j|j \neq i} \sim q_\alpha(\mathbf{t}_2)}}} \left[\log \frac{f(\mathbf{t}_{2,i}, \mathbf{t}_{1,i})}{\sum_j f(\mathbf{t}_{2,j}, \mathbf{t}_{1,i})} \right] \quad (2)$$

where $f(\mathbf{t}_{2,j}, \mathbf{t}_{1,i}) = \exp(\mathbf{g}(\mathbf{t}_{2,j}) \cdot \mathbf{g}(\mathbf{t}_{1,i})/\tau)$ and $\mathbf{g}(\cdot)$ represents a multimodal feature encoder and τ is a temperature parameter. We provide an example for the TupleInfoNCE objective in Figure 3. The hyper-parameters α and β can be optimized to allow flexible control over the contribution of different modalities as introduced in the next section.

Connection with Mutual Information estimation To better understand why $\mathcal{L}_{\text{TNCE}}^{\alpha\beta}$ is more suited for multimodal fusion than \mathcal{L}_{NCE} , we provide a theoretical analysis from the information theory perspective. As we mentioned in

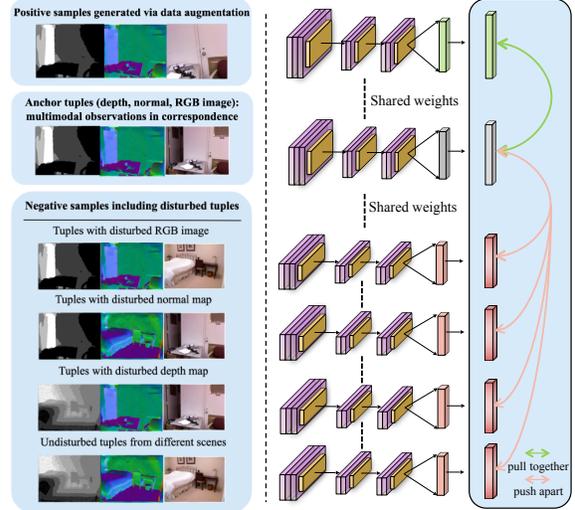


Figure 3. An example of the TupleInfoNCE objective for RGB, depth and normal map fusion.

Section 3, minimizing \mathcal{L}_{NCE} is equivalent to maximizing a lower bound of $I(\mathbf{t}_2; \mathbf{t}_1)$, which could lead to weak modalities and the modality synergy being ignored. Minimizing $\mathcal{L}_{\text{TNCE}}^{\alpha\beta}$, instead, is equivalent to maximizing a lower bound of $I(\mathbf{t}_2; \mathbf{t}_1|\beta) + \sum_{k=1}^K \alpha_k I(\mathbf{v}_2^k; \bar{\mathbf{v}}_2^k)$ (please see supplementary material for a proof). As is shown in Figure 2, $I(\mathbf{v}_2^k; \bar{\mathbf{v}}_2^k)$ puts more emphasis on the information shared across modalities to encourage modality synergy and to avoid weak modalities being ignored. The ratio of k -disturbed negative samples α_k plays the role of balancing $I(\mathbf{v}_2^k; \bar{\mathbf{v}}_2^k)$ and $I(\mathbf{t}_2; \mathbf{t}_1|\beta)$. And the data augmentation parameters β directly influence $I(\mathbf{t}_2; \mathbf{t}_1|\beta)$ and further balance the information contribution of each modality.

4.3. Sample Optimization

The hyper-parameters α and β designed for tuple disturbing and augmentation play a key role in the TupleInfoNCE objective design. Each set of α and β will correspond to one specific objective and fully optimizing $\mathcal{L}_{\text{TNCE}}^{\alpha\beta}$ will result in a multimodal feature encoder $\mathbf{g}^{\alpha\beta}$. Manually setting these hyper-parameters is not reliable, motivating us to explore ways to optimize these hyper-parameters. There are mainly two challenges to be addressed. The first is the evaluation challenge: we need a way to evaluate the quality of the multimodal feature encoder $\mathbf{g}^{\alpha\beta}$ in an unsupervised manner since most existing works have demonstrated that InfoNCE loss itself is not a good evaluator [32, 22]. The second is the optimization challenge: we need an efficient optimization strategy to avoid exhaustively examining different hyper-parameters and training the whole network from scratch repeatedly. We will explain how we handle these challenges to optimize the ratio α of different types of negative samples in Section 4.3.1, and the hyper-parameter β of augmented positive samples in Section 4.3.2.

4.3.1 Optimizing negative samples

To evaluate the modality fusion quality in the learned representations unsupervisedly, we propose to use crossmodal discrimination as a surrogate task. To efficiently optimize α , we adopt a bilevel optimization scheme alternating between optimizing α and optimizing the main $\mathcal{L}_{\text{TNCE}}^{\alpha\beta}$ objective with a fixed α . We elaborate on these designs below.

Crossmodal discrimination TupleInfoNCE differs from the naive InfoNCE in that it emphasizes more on each modality \mathbf{v}^k as well as its mutual information $I(\mathbf{v}^k; \bar{\mathbf{v}}^k)$ with the rest modalities $\bar{\mathbf{v}}^k$. In order to learn a good representation that properly covers $I(\mathbf{v}^k; \bar{\mathbf{v}}^k)$, we propose a novel surrogate task, *crossmodal discrimination*, which looks for the corresponding $\bar{\mathbf{v}}^k$ only by examining \mathbf{v}^k in a holdout validation set. Mathematically, we first generate a validation set $\{\mathbf{t}_m\}_{m=1}^M$ by drawing M random tuples $\mathbf{t}_m = (\mathbf{v}_m^1, \mathbf{v}_m^2, \dots, \mathbf{v}_m^K) \sim p(\mathbf{t})$. For each modality \mathbf{v}_m^k , its augmented version is represented as $\mathbf{v}_m^{k'} \sim p_{\zeta_k}(\mathbf{v}_m^k | \mathbf{v}_m^k)$ following a data augmentation strategy parameterized by ζ_k . Then the crossmodal discrimination task is defined as, given any $\mathbf{v}_n^{k'}$ sampled from the augmented validation set $\{\mathbf{v}_m^{k'}\}_{m=1}^M$, finding its corresponding rest modalities $\bar{\mathbf{v}}_n^k$ in the set $\{\bar{\mathbf{v}}_m^k\}_{m=1}^M$. To solve this surrogate task, for any $\mathbf{v}_n^{k'}$ sampled from the augmented validation set $\{\mathbf{v}_m^{k'}\}_{m=1}^M$, we first compute its probability that corresponds to $\bar{\mathbf{v}}_l^k$ as,

$$p_{nl}^k(\mathbf{g}^{\alpha\beta}) = \frac{\exp(\mathbf{g}^{\alpha\beta}(\mathbf{v}_n^{k'}) \cdot \mathbf{g}^{\alpha\beta}(\bar{\mathbf{v}}_l^k) / \tau)}{\sum_{m=1}^M \exp(\mathbf{g}^{\alpha\beta}(\mathbf{v}_n^{k'}) \cdot \mathbf{g}^{\alpha\beta}(\bar{\mathbf{v}}_m^k) / \tau)} \quad (3)$$

where $\mathbf{g}^{\alpha\beta}(\cdot)$ represents our optimal multimodal feature encoder trained via optimizing $\mathcal{L}_{\text{TNCE}}^{\alpha\beta}$ and τ is a temperature parameter. Then the crossmodal discrimination accuracy for the k -th modality can be computed as

$$\mathcal{A}^k(\mathbf{g}^{\alpha\beta}) = \sum_{n=1}^M \mathbb{1}(n = \arg \max_l p_{nl}^k(\mathbf{g}^{\alpha\beta})) / M \quad (4)$$

where $\mathbb{1}(\cdot)$ is an indicator function. $\mathcal{A}^k(\mathbf{g}^{\alpha\beta})$ roughly measures how much $I(\mathbf{v}^k; \bar{\mathbf{v}}^k)$ the encoder $\mathbf{g}^{\alpha\beta}$ has captured and provides cues regarding how we should adjust α_k in the negative samples. We can then leverage the crossmodal discrimination accuracy to optimize α through maximizing the following reward:

$$\mathcal{R}(\alpha) = \sum_{k=1}^K \mathcal{A}^k(\mathbf{g}^{\alpha\beta}) \quad (5)$$

which properly balances the contribution of different modalities and has a high correlation with downstream semantic inference tasks as shown in Section 5.4. Notice to handle missing modalities in the crossmodal discrimination task, we adopt a *dropout training* strategy as introduced in the supplemental material.

Bilevel optimization Now we describe how to efficiently optimize $\mathcal{R}(\alpha)$ with one-pass network training. We write

our optimization problem as below:

$$\begin{aligned} \text{maximize } \mathcal{R}(\alpha) &= \sum_{k=1}^K \mathcal{A}^k(\mathbf{g}^{\alpha\beta}) \\ \text{s.t. } \mathbf{g}^{\alpha\beta} &= \arg \min_{\mathbf{g}} \mathcal{L}_{\text{TNCE}}^{\alpha\beta}(\mathbf{g}) \end{aligned} \quad (6)$$

This is a standard bilevel optimization problem. Inspired by [16], we adopt a hyper-parameter optimization strategy which alternatively optimizes α and \mathbf{g} in a single training pass. Specifically, we relax the constraint that $\sum_{k=0}^K \alpha_k = 1$ during the optimization and use an independent multivariate Gaussian $\mathcal{N}(\mu_0, \sigma I)$ to initialize the distribution of α . At each training epoch t , we sample B hyper-parameters $\{\alpha_1, \dots, \alpha_B\}$ from distribution $\mathcal{N}(\mu_t, \sigma I)$ and train our current feature encoder \mathbf{g}_t separately to generate B new encoders $\{\mathbf{g}_{t+1}^1, \dots, \mathbf{g}_{t+1}^B\}$. We evaluate the reward for each of these encoders on the validation set and update the distribution of α using REINFORCE [35] as below:

$$\mu_{t+1} = \mu_t + \eta \frac{1}{B} \sum_{i=1}^B R(\alpha_i) \nabla_{\alpha} \log(p(\alpha_i; \mu, \sigma)) \quad (7)$$

where $p(\alpha_i; \mu, \sigma)$ represents the PDF of the Gaussian distribution. We then pick up the encoder with the highest reward as our \mathbf{g}_{t+1} and continue with the next epoch. We repeat the above process until convergence.

4.3.2 Optimizing positive samples

Similar to optimizing α , a reward function is required to evaluate our feature encoder $\mathbf{g}^{\alpha\beta}$ in an unsupervised manner with respect to β . A straightforward approach is to adopt the total crossmodal discrimination accuracy defined in Equation 5. Through experiments, we observe two phenomena making this simple adaptation fail to optimize β effectively. We use β and ζ to represent the data augmentation parameters for training and validation respectively, and they do not have to be the same. 1). If we manually set ζ to be fixed, the optimal β maximizing the total accuracy highly correlates with ζ and fails to generate truly good positive samples. 2). If we set ζ to be the same as β and optimize them together, we usually achieve the best total accuracy when no data augmentation is applied, though it has been shown a certain level of data augmentation is important for contrastive learning [5, 32]. Therefore a better reward function is required for β optimization.

We re-write our total crossmodal discrimination accuracy as $\sum_{k=1}^K \mathcal{A}^k(\mathbf{g}^{\alpha\beta}, \zeta)$ to reflect the influence from ζ . Instead of manually setting ζ which produces a chicken-and-egg problem for hyper-parameter optimization, we set $\zeta = \beta$ and only optimize β . We follow the conclusion in [32] and aim to use strong data augmentations, which reduces the information contribution by each modality but make the contributed information more robust to nuanced input noises.

Algorithm 1: Sample Optimization

Input: Initialized multimodal feature encoder \mathbf{g}_0 , initialized distribution $(\mu_0^\alpha, \sigma^\alpha)$ and $(\mu_0^\beta, \sigma^\beta)$, total training epochs T , distribution learning rate η

Output: Final multimodal feature encoder $\mathbf{g}_T^{\alpha^* \beta^*}$

for $t = 1$ **to** T **do**

if t is even **then**

 Sample B **sampling ratio** hyper-parameters

$\{\alpha_i\}_{i=1}^B$ via distribution $\mathcal{N}(\mu_t^\alpha, \sigma^\alpha I)$;

 Train \mathbf{g}_t for one epoch separately with each α_i and

 get $\{\mathbf{g}_{t+1}^i\}_{i=1}^B$;

 Calculate rewards $\{\mathcal{R}(\alpha_i)\}_{i=1}^B$ using Equation 5;

 Decide the best model $i = \arg \max_j \mathcal{R}(\alpha_j)$;

 Update μ_{t+1}^α using Equation 7;

 Update $\mathbf{g}_{t+1} = \mathbf{g}_{t+1}^i$;

else if t is odd **then**

 Sample B **data augmentation** hyper-parameters

$\{\beta_i\}_{i=1}^B$ via distribution $\mathcal{N}(\mu_t^\beta, \sigma^\beta I)$;

 Train \mathbf{g}_t for one epoch separately with each β_i and

 get $\{\mathbf{g}_{t+1}^i\}_{i=1}^B$;

 Calculate rewards $\{\mathcal{R}(\beta_i)\}_{i=1}^B$ using Equation 8;

 Decide the best model $i = \arg \max_j \mathcal{R}(\beta_j)$;

 Update μ_{t+1}^β using Equation 7;

 Update $\mathbf{g}_{t+1} = \mathbf{g}_{t+1}^i$;

end if

end for

return \mathbf{g}_T

We observe that the total accuracy will decrease as we use stronger augmentations, and minimizing $\sum_{k=1}^K \mathcal{A}^k(\mathbf{g}^{\alpha\beta}, \beta)$ with respect to β will effectively increase the augmentation magnitude. However, as discussed in [32], we should not increase the data augmentation without any constraints and there is a sweet spot going beyond which a larger data augmentation could harm the representation learning. We find $\|\beta - \zeta^*(\beta)\|^2$ providing cues for identifying the sweet spot, where $\zeta^*(\beta) = \arg \max_{\zeta} \sum_{k=1}^K \mathcal{A}^k(\mathbf{g}^{\alpha\beta}, \zeta)$ represents the best ζ maximizing the total crossmodal discrimination accuracy $\sum_{k=1}^K \mathcal{A}^k$ for a feature encoder trained with β . When β is weak, we empirically discover that $\zeta^*(\beta)$ is very close to β ; when β is too strong, smaller augmentation parameters on the validation set will lead to higher total accuracy, therefore leading to a large difference between β and $\zeta^*(\beta)$. We provide empirical studies supporting these findings in Section 5.4. Motivated by the above observations, we design our reward function as:

$$\mathcal{R}(\beta) = 1 - \sum_{k=1}^K \frac{\mathcal{A}^k(\mathbf{g}^{\alpha\beta}, \beta)}{K} - \lambda \frac{\|\beta - \zeta^*(\beta)\|^2}{\|\beta^{\max}\|^2} \quad (8)$$

where λ is a balancing parameter and β^{\max} denotes a pre-defined augmentation parameter upper bound used for the normalization purpose.

$\mathcal{R}(\beta)$ can be optimized in the same way as how $\mathcal{R}(\alpha)$ is optimized, and we alternate between optimizing β and \mathbf{g} in a single training pass. We further combine the optimization of $\mathcal{R}(\alpha)$, $\mathcal{R}(\beta)$, and the multimodal encoder \mathbf{g} in Algorithm 1, where we update α when the epoch number is even and update β otherwise.

5. Experiment

In this section, we evaluate our method by transfer learning, i.e., fine-tuning on downstream tasks and datasets. Specifically, we first pretrain our backbone on each dataset without any additional data using the proposed TupleInfoNCE. Then we use the pre-trained weights as initialization and further refine them for target downstream tasks. In this case, good features could directly lead to performance gains in downstream tasks.

We present results for three popular multi-modality tasks: semantic segmentation on NYUv2 [29], 3D object detection on SUN RGB-D [30], and sentiment analysis on MOSEI [39] and MOSI [38] in Section 5.1, 5.2 and 5.3 respectively. In Section 5.4, extensive ablation studies, analysis and visualization are provided to justify design choices of our system.

5.1. NYUv2 Semantic Segmentation

Setup. We first conduct experiments on NYUv2 [29] to see whether our method can help multimodal semantic scene understanding. NYUv2 contains 1,449 indoor RGB-D images, of which 795 are used for training and 654 for testing. We use three modalities in this task: RGB, depth, and normal map. The data augmentation strategies we adopted include random cropping, rotation, and color jittering. We use ESANet [28], an efficient ResNet-based encoder, as our backbone. We use the common 40-class label setting and mean IoU(mIoU) as the evaluation metric.

We compare our method with the train-from-scratch baseline as well as the latest self-supervised multimodal representation learning methods including CMC [31], MMV FAC [1] and MISA [11], which are all based upon crossmodal embedding. In addition, we include an InfoNCE [24] baseline where we directly contrast multimodal input tuples without tuple disturbing and sample optimization. We also include supervised pretraining [28] methods for completeness.

Results. Table 1 shows that the previous best performing method MISA [11] improves the segmentation mIoU by 3.3% over the train-from-scratch baseline. When using InfoNCE [24], the improvement drops to 2.0%. Our method achieves 8.0% improvement over the train-from-scratch baseline. The improvement from 40.1% to 48.1% confirms that we can produce better-fused representations to boost the segmentation performance on RGB-D scenes. Notably, our proposed TupleNCE, though only pretrained on NYUv2 self-supervisedly, is only ~3% lower than supervised pretraining methods.

Table 1. Semantic Segmentation results on NYUv2.

Methods	mIoU
Train from scratch	40.1
Supervised pretrain on Imagenet	50.3
Supervised pretrain on Scenenet	51.6
CMC	41.9
MMV FAC	42.5
MISA	43.4
InfoNCE	42.1
Ours	48.1

5.2. SUN RGB-D 3D Object Detection

Setup. Our second experiment investigates how Tuple-InfoNCE can be used for 3D object detection in the SUN RGB-D dataset [30]. SUN RGB-D contains a training set with ~5K single-view RGB-D scans and a test set with ~5K scans. The scans are annotated with amodal 3D-oriented bounding boxes for objects from 37 categories. We use three modalities in this experiment: 3D point cloud, RGB color and height. Data augmentation used here is rotation for point cloud, jittering for RGB color, and random noise for height. We use VoteNet [25] as our backbone, which leverages PointNet++ [26] to process depth point cloud and supports appending RGB or height information as additional inputs. We compare our method with baseline methods including InfoNCE [24], CMC [31], and MISA [11]. We use mAP@0.25 as our evaluation metric.

Table 2. 3D Object Detection results on SUN RGB-D.

Methods	mAP@0.25
Train from scratch	56.3
InfoNCE	56.8
CMC	56.5
MISA	56.7
Ours	58.0

Results. Table 2 shows the object detection results. We find that previous self-supervised methods seem to struggle with 3D tasks: CMC and MISA achieve very limited improvement over the baseline trained from scratch. The improvement of InfoNCE [24] is also very marginal (0.5%), presumably because overemphasizing the modality-specific information from strong modalities might sacrifice the weak modalities as well as the modality synergy during learning. In contrast, TupleInfoNCE achieves 1.7% mAP improvement over the baseline trained from scratch, which more than triples the improvement InfoNCE achieved. The comparison between our method and InfoNCE directly validates the efficacy of the proposed TupleInfoNCE objective and sample optimization mechanism.

5.3. Multimodal Sentiment Analysis

Setup. Our third experiment investigates multimodal sentiment analysis with the MOSI [38] and MOSEI [39] datasets, both providing word-aligned multimodal signals (language, visual and acoustic) for each utterance. MOSI contains 2198 subjective utterance-video segments. The utterances are manually annotated with a continuous opinion score between [-3,3], where -3/+3 represents strongly negative/positive sentiments. MOSEI is an improvement over MOSI with a higher number of utterances, greater variety in samples, speakers, and topics. Following the recent state-of-the-art multimodal self-supervised representation learning method MISA [11], we use features pre-extracted from the original raw data, which does not permit an intuitive way for data augmentation. Therefore we only optimize negative samples in this experiment. We use the same backbone as MISA [11] to make a fair comparison. We use binary accuracy (Acc-2), 7-class accuracy (Acc-7), and F-Score as our evaluation metrics.

Results. As shown in Table 3 and 4, our method consistently outperforms previous methods on these very challenging and competitive datasets – e.g., compared with the previous best performing method MISA, the Acc-7 goes up from 42.3 to 43.3 on MOSI, and from 52.2 to 52.7 on MOSEI. As these two approaches share the same network backbone and only differ in their strategy to learn the fused representation, the improvement provides strong evidence for the effectiveness of our method.

Table 3. Multimodal sentiment analysis results on MOSI.

Methods	Acc-2	Acc-7	F-Score
Train from scratch	83.0	40.0	82.8
CMC	83.3	39.5	83.0
MMV FAC	83.5	41.5	83.4
MISA	83.4	42.3	83.6
InfoNCE	83.1	40.5	82.8
Ours	83.6	43.3	83.8

Table 4. Multimodal sentiment analysis results on MOSEI.

Methods	Acc-2	Acc-7	F-Score
Train from scratch	82.5	51.8	82.3
CMC	83.3	50.8	84.1
MMV FAC	85.1	52.0	85.0
MISA	85.5	52.2	85.3
InfoNCE	83.5	52.0	83.4
Ours	86.1	52.7	86.0

5.4. Further Analysis and Discussions

Efficacy of sample optimization We run ablation studies with and without sample optimization to quantify its efficacy. We find that uniformly setting α_k without optimizing neg-

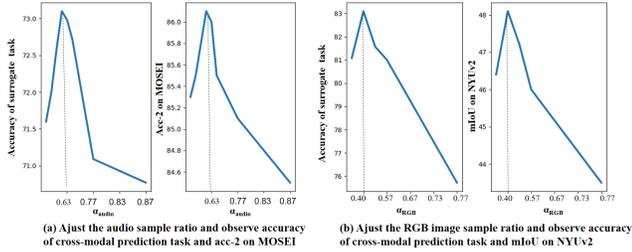


Figure 4. Correlations between the total crossmodal discrimination accuracy and the downstream task performance.

ative samples results in a 1.7% mIoU drop on the NYUv2 semantic segmentation task, 0.5 mAP drop on the SUN RGB-D 3D object detection task, 0.6 Acc-7 drop on MOSI, and 0.4 Acc-7 drop on MOSEI. Manually designing data augmentation strategies without optimizing positive samples as in [31] results in a 1.1 mIoU drop on NYUv2 and 0.6 mAP drop on SUN RGB-D. We also examine the optimized negative sampling strategy as well as the data augmentation strategy. On the NYUv2 dataset, we find the best performing negative sampling ratio among RGB, depth and normal is roughly 2 : 1 : 1, showing that RGB is emphasized more in the fused representations. As for the data augmentation strategy, though we use the same types of data augmentations for all the three modalities on NYUv2, the optimal augmentation parameters vary from modality to modality. Considering image rotation with the hyper-parameter representing the rotation angle, we found that 40 degrees is the best hyper-parameter for RGB images, while 10 degrees is the best for depth and normal maps. Please refer to the supplementary material for more analysis regarding SUN RGB-D, MOSI, and MOSEI.

Reward design for negative sample optimization We introduce crossmodal discrimination as a surrogate task for negative sample optimization in Section 4.3.1 and argue that the total crossmodal discrimination accuracy $\mathcal{R}(\alpha)$ in Equation 5 is a good reward function. We provide our empirical verification here. We vary the ratio α_k of type- k negative samples while keeping the relative ratio of the rest types unchanged. We train the whole network through with the fixed negative sampling ratio and evaluate both $\mathcal{R}(\alpha)$ and the performance of the downstream task. As is shown in Figure 4, adjusting the proportion of different types of negative samples will influence the accuracy $\mathcal{R}(\alpha)$ of the surrogate task, which has a high correlation with downstream tasks. Too low and too high proportion for one type of negative samples both lead to low $\mathcal{R}(\alpha)$. There is a sweet spot corresponding to the best $\mathcal{R}(\alpha)$. Experiments show this sweet spot also corresponds to the best performance on downstream tasks.

Reward design for positive sample optimization Our reward function in Equation 8 for positive sample optimization is motivated by two observations: 1). minimizing total crossmodal discrimination accuracy $\sum_{k=1}^K A^k(\mathbf{g}^{\alpha\beta}, \beta)$ with respect to β will increase the augmentation magnitude; 2).

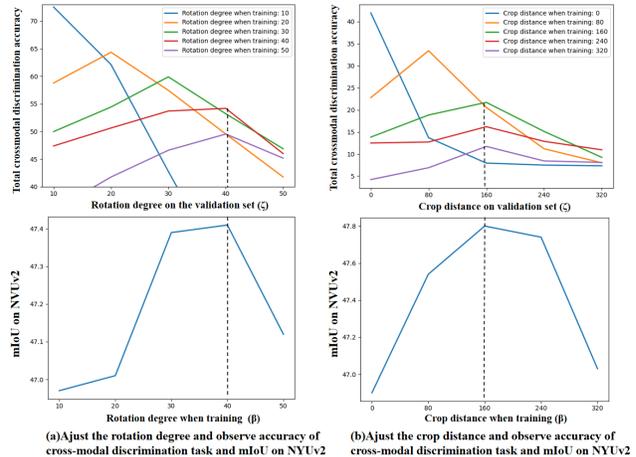


Figure 5. Empirical study justifying the reward design for positive sample optimization. In the first row we show the total crossmodal discrimination accuracy on the validation set while varying the augmentation parameter ζ and different curves are obtained with different train time data augmentation parameters β . The second row shows how the performance of downstream tasks vary while changing β .

$\|\beta - \zeta^*(\beta)\|^2$ provides cues for identifying the sweet spot beyond which larger augmentation will harm representation learning. We provide empirical studies to verify these observations in Figure 5. We train networks from beginning and end with different β to evaluate how the total crossmodal discrimination accuracy change while varying the data augmentation parameters ζ on the validation set. We also evaluate how the performance of downstream tasks varies while changing the training time data augmentation parameters β . We experiment with two types of data augmentation - image rotation and image crop, and obtain consistent observations. $\sum_{k=1}^K A^k(\mathbf{g}^{\alpha\beta}, \beta)$ indeed drops while increasing β . Moreover, $\zeta^*(\beta)$ corresponds to the peak of each curve in the first row and it is very close to β when β is small. Once β goes beyond a sweet spot, which gives the best performance on downstream tasks, $\zeta^*(\beta)$ no longer tracks the value of β and $\|\beta - \zeta^*(\beta)\|^2$ will give a penalty for further increasing β . In practice, we find our reward function powerful enough for identifying the best training time data augmentation parameters.

Robustness to uninformative modality TupleInfoNCE emphasizes the modality which is easy to be ignored. An obvious question is whether it is robust to uninformative modalities. We conduct experiments on MOSEI multimodal sentiment analysis task and add an uninformative modality named timestamp which denotes the relative time in a sequence. Results show using these four modalities, we achieve 52.6 Acc-7, which is only 0.1% lower than before. The final negative sample ratio among the four modalities is roughly 3(text): 3(video): 4(audio): 1(timestamp), showing our method successfully identifies that “timestamp” is not something worthy of much emphasis.

6. Conclusion

This paper proposes a new objective for representation learning of multimodal data using contrastive learning, TupleInfoNCE. The key idea is to contrast multimodal anchor tuples with challenging negative samples containing disturbed modalities and better positive samples obtained through an optimizable data augmentation process. We provide a theoretical basis for why TupleInfoNCE works, an algorithm for optimizing TupleInfoNCE with a self-supervised approach to select the contrastive samples, and results of experiments showing ablations and state-of-the-art performance on a wide range of multimodal fusion benchmarks.

Acknowledgement

This work was supported by the Key-Area Research and Development Program of Guangdong Province (2019B121204008) and the Center on Frontiers of Computing Studies (7100602567).

References

- [1] Jean-Baptiste Alayrac, Adrià Recasens, Rosalia Schneider, Relja Arandjelović, Jason Ramapuram, Jeffrey De Fauw, Lucas Smaira, Sander Dieleman, and Andrew Zisserman. Self-supervised multimodal versatile networks. *arXiv preprint arXiv:2006.16228*, 2020. 1, 2, 6
- [2] Relja Arandjelovic and Andrew Zisserman. Objects that sound. In *Proceedings of the European conference on computer vision (ECCV)*, pages 435–451, 2018. 2
- [3] Bowen Baker, Otkrist Gupta, Nikhil Naik, and Ramesh Raskar. Designing neural network architectures using reinforcement learning. *arXiv preprint arXiv:1611.02167*, 2016. 2
- [4] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *arXiv preprint arXiv:2006.09882*, 2020. 2
- [5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 1, 2, 5
- [6] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 2
- [7] Ching-Yao Chuang, Joshua Robinson, Lin Yen-Chen, Antonio Torralba, and Stefanie Jegelka. Debaised contrastive learning. *arXiv preprint arXiv:2007.00224*, 2020. 2
- [8] Joon Son Chung and Andrew Zisserman. Out of time: automated lip sync in the wild. In *Asian conference on computer vision*, pages 251–263. Springer, 2016. 2
- [9] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*, 2018. 2
- [10] Basura Fernando, Hakan Bilen, Efstratios Gavves, and Stephen Gould. Self-supervised video representation learning with odd-one-out networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3636–3645, 2017. 2
- [11] Devamanyu Hazarika, Roger Zimmermann, and Soujanya Poria. Misa: Modality-invariant and-specific representations for multimodal sentiment analysis. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 1122–1131, 2020. 2, 6, 7
- [12] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020. 1, 2
- [13] Chih-Hui Ho and Nuno Vasconcelos. Contrastive learning with adversarial examples. *arXiv preprint arXiv:2010.12050*, 2020. 2
- [14] Tomas Jakab, Ankush Gupta, Hakan Bilen, and Andrea Vedaldi. Unsupervised learning of object landmarks through conditional image generation. *arXiv preprint arXiv:1806.07823*, 2018. 2
- [15] Yannis Kalantidis, Mert Bulent Sariyildiz, Noe Pion, Philippe Weinzaepfel, and Diane Larlus. Hard negative mixing for contrastive learning. *arXiv preprint arXiv:2010.01028*, 2020. 2
- [16] Chuming Li, Xin Yuan, Chen Lin, Minghao Guo, Wei Wu, Junjie Yan, and Wanli Ouyang. Am-lfs: Automl for loss function search. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8410–8419, 2019. 2, 5
- [17] Junnan Li, Pan Zhou, Caiming Xiong, Richard Socher, and Steven CH Hoi. Prototypical contrastive learning of unsupervised representations. *arXiv preprint arXiv:2005.04966*, 2020. 2
- [18] Sungbin Lim, Ildoo Kim, Taesup Kim, Chiheon Kim, and Sungwoong Kim. Fast autoaugment. In *NeurIPS*, 2019. 2
- [19] Chenxi Liu, Piotr Dollár, Kaiming He, Ross Girshick, Alan Yuille, and Saining Xie. Are labels necessary for neural architecture search? In *European Conference on Computer Vision*, pages 798–813. Springer, 2020. 3
- [20] Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. In *Proceedings of the European conference on computer vision (ECCV)*, pages 19–34, 2018. 2
- [21] Yunze Liu, Li Yi, Shanghang Zhang, Qingnan Fan, Thomas Funkhouser, and Hao Dong. P4contrast: Contrastive learning with pairs of point-pixel pairs for rgb-d scene understanding. *arXiv preprint arXiv:2012.13089*, 2020. 1, 2, 3
- [22] Sean Metzger, Aravind Srinivas, Trevor Darrell, and Kurt Keutzer. Evaluating self-supervised pretraining without using labels. *arXiv preprint arXiv:2009.07724*, 2020. 3, 4
- [23] Ishan Misra, C Lawrence Zitnick, and Martial Hebert. Shuffle and learn: unsupervised learning using temporal order verification. In *European Conference on Computer Vision*, pages 527–544. Springer, 2016. 2
- [24] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 1, 2, 3, 6, 7

- [25] Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. Deep hough voting for 3d object detection in point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9277–9286, 2019. 7
- [26] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *arXiv preprint arXiv:1706.02413*, 2017. 7
- [27] Joshua Robinson, Ching-Yao Chuang, Suvrit Sra, and Stefanie Jegelka. Contrastive learning with hard negative samples. *arXiv preprint arXiv:2010.04592*, 2020. 2
- [28] Daniel Seichter, Mona Köhler, Benjamin Lewandowski, Tim Wengefeld, and Horst-Michael Gross. Efficient rgb-d semantic segmentation for indoor scene analysis. *arXiv preprint arXiv:2011.06961*, 2020. 6
- [29] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *European conference on computer vision*, pages 746–760. Springer, 2012. 2, 6
- [30] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 567–576, 2015. 2, 6, 7
- [31] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019. 1, 2, 6, 7, 8, 11
- [32] Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What makes for good views for contrastive learning. *arXiv preprint arXiv:2005.10243*, 2020. 2, 3, 4, 5, 6, 11
- [33] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103, 2008. 2
- [34] Donglai Wei, Joseph J Lim, Andrew Zisserman, and William T Freeman. Learning and using the arrow of time. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8052–8060, 2018. 2
- [35] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992. 5
- [36] Mike Wu, Chengxu Zhuang, Milan Mosse, Daniel Yamins, and Noah Goodman. On mutual information in contrastive learning for visual representations. *arXiv preprint arXiv:2005.13149*, 2020. 2
- [37] Saining Xie, Jiatao Gu, Demi Guo, Charles R Qi, Leonidas Guibas, and Or Litany. Pointcontrast: Unsupervised pre-training for 3d point cloud understanding. In *European Conference on Computer Vision*, pages 574–591. Springer, 2020. 2, 3
- [38] Amir Zadeh, Rowan Zellers, Eli Pincus, and Louis-Philippe Morency. Multimodal sentiment intensity analysis in videos: Facial gestures and verbal messages. *IEEE Intelligent Systems*, 31(6):82–88, 2016. 2, 6, 7
- [39] AmirAli Bagher Zadeh, Paul Pu Liang, Soujanya Poria, Erik Cambria, and Louis-Philippe Morency. Multimodal language analysis in the wild: Cmu-mosei dataset and interpretable dynamic fusion graph. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2236–2246, 2018. 2, 6, 7
- [40] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8697–8710, 2018. 2

Appendix

A. Does representation quality improve as number of views increases?

To measure the quality of the learned representation, we consider the task of semantic segmentation on NYUv2. In the 1 modality setting case, TupleInfoNCE using RGB modality coincides with InfoNCE. In 2-3 modalities cases, we sequentially add depth and normal modalities.

Modality	mIoU
RGB	40.8
RGB + depth	44.1
RGB + normal	44.5
depth + normal	46.3
RGB + depth + normal	48.1

Table 5. We show the mean Intersection over Union (mIoU) for the NYU-Depth-V2 dataset, as TupleInfoNCE is trained with increasingly more modalities from 1 to 3. The performance steadily improves as new modalities are added.

As shown in Tab. 6, We see that the performance steadily improves as new modalities are added. This finding is consistent with that of CMC [31] who learn a cross-modal embedding space.

B. Comparisons to InfoMin

Infomin [32] proposes to reduce the mutual information between different views while retaining the task-relevant information as much as possible. But how to explore task-relevant information without labels is a very challenging problem. Infomin utilizes an adversarial training strategy to search for good views in a weakly-supervised manner. When no labels are available, which is the case in a truly self-supervised representation learning setting, the effectiveness of Infomin is greatly reduced. Taking NYUv2 semantic segmentation task as an example, in a truly self-supervised representation learning setting, the optimal rotation parameter Infomin finds is 70 degrees, which leads to 46.91 mIoU after downstream fine-tuning, while the optimal rotation parameter TupleInfoNCE finds is 40 degree, corresponding to 47.41 mIoU.

C. Use TupleInfoNCE in four modalities

We conduct another experiment on NYUv2 to see whether our method can help multi-modal semantic scene understanding in the case of 4 modalities. Following CMC, modalities we used are L, ab, depth, normal. We follow the 3 modalities setting and use random cropping, rotation, and jittering as the augmentation strategy, thus learning a fused representation by contrasting tuples that contain 4 modalities. Using 4 modalities, the best mIoU we obtain is 48.6. This

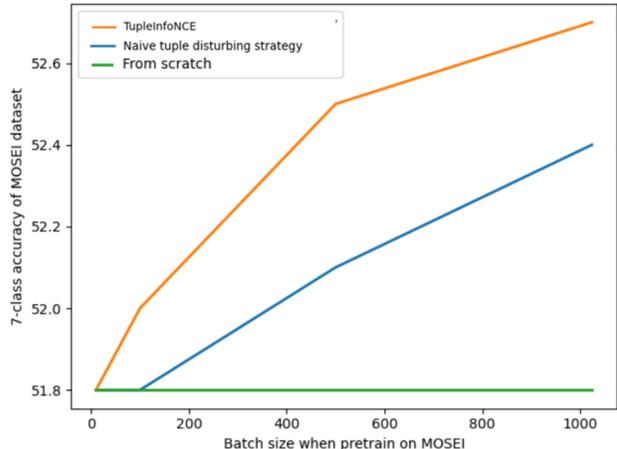


Figure 6. Sample efficiency of TupleInfoNCE

is 0.5 higher than our 3 modalities baseline, showing that our method has strong generalization ability. These again validate that our TupleInfoNCE which contrasts multi-modal input tuples can produce better-fused representations.

D. Sampling efficiency of TupleInfoNCE

TupleInfoNCE disturbs each modality separately to generate k disturbed negative samples. However, The most direct method is the naive sampling strategy which disturbs all modalities simultaneously to generate 1 disturbed negative sample. We conduct another experiments to compare our method with the naive sampling strategy from the perspective of efficiency. Figure 6 shows that our method is more efficient than naive tuple disturbing strategy. Our method with a batch size of 512 already outperforms naive tuple disturbing strategy with a batch size of 1024.

E. Proof for the MI lower bound of TupleInfoNCE

The TupleInfoNCE loss is essentially a categorical cross-entropy classifying positive tuples $\mathbf{t}_{2,i} \sim p_{\beta}(\mathbf{t}_2|\mathbf{t}_{1,i})$ from $N-1$ negatives views $\mathbf{t}_{2,j} \sim q_{\alpha}(\mathbf{t}_2)$. We use $p_{\alpha\beta}(d=i|\mathbf{t}_{1,i})$ to denote the optimal probability for this loss where $[d=i]$ is an indicator showing that $\mathbf{t}_{2,i}$ is the positive view. Then we can derive $p_{\alpha\beta}(d=i|\mathbf{t}_{1,i})$ as follows:

$$\begin{aligned}
 p_{\alpha\beta}(d=i|\mathbf{t}_{1,i}) &= \frac{p_{\beta}(\mathbf{t}_{2,i}|\mathbf{t}_{1,i}) \prod_{l \neq i} q_{\alpha}(\mathbf{t}_{2,l})}{\sum_{j=1}^N p_{\beta}(\mathbf{t}_{2,j}|\mathbf{t}_{1,i}) \prod_{l \neq j} q_{\alpha}(\mathbf{t}_{2,l})} \\
 &= \frac{p_{\beta}(\mathbf{t}_{2,i}|\mathbf{t}_{1,i})}{q_{\alpha}(\mathbf{t}_{2,i})} \\
 &= \frac{p_{\beta}(\mathbf{t}_{2,i}|\mathbf{t}_{1,i})}{\sum_{j=1}^N \frac{p_{\beta}(\mathbf{t}_{2,j}|\mathbf{t}_{1,i})}{q_{\alpha}(\mathbf{t}_{2,j})}}
 \end{aligned}$$

It can be seen from above that the optimal value for $f(\mathbf{t}_{2,j}, \mathbf{t}_{1,i})$ in $\mathcal{L}_{\text{TNCE}}(\alpha, \beta)$ is proportional to $\frac{p_{\beta}(\mathbf{t}_{2,j}|\mathbf{t}_{1,i})}{q_{\alpha}(\mathbf{t}_{2,j})}$. Insert this density ratio back to $\mathcal{L}_{\text{TNCE}}(\alpha, \beta)$ we get:

$$\begin{aligned}
& \mathcal{L}_{\text{TNCE}}^{\text{OPT}}(\alpha, \beta) \\
&= -\mathbb{E} \log \left[\frac{\frac{p_{\beta}(\mathbf{t}_{2,i}|\mathbf{t}_{1,i})}{q_{\alpha}(\mathbf{t}_{2,i})}}{\frac{p_{\beta}(\mathbf{t}_{2,i}|\mathbf{t}_{1,i})}{q_{\alpha}(\mathbf{t}_{2,i})} + \sum_{j \neq i} \frac{p_{\beta}(\mathbf{t}_{2,j}|\mathbf{t}_{1,i})}{q_{\alpha}(\mathbf{t}_{2,j})}} \right] \\
&= \mathbb{E} \log \left[1 + \frac{q_{\alpha}(\mathbf{t}_{2,i})}{p_{\beta}(\mathbf{t}_{2,i}|\mathbf{t}_{1,i})} \sum_{j \neq i} \frac{p_{\beta}(\mathbf{t}_{2,j}|\mathbf{t}_{1,i})}{q_{\alpha}(\mathbf{t}_{2,j})} \right] \\
&\approx \mathbb{E} \log \left[1 + \frac{q_{\alpha}(\mathbf{t}_{2,i})}{p_{\beta}(\mathbf{t}_{2,i}|\mathbf{t}_{1,i})} (N-1) \mathbb{E}_{\mathbf{t}_{2,j}} \frac{p_{\beta}(\mathbf{t}_{2,j}|\mathbf{t}_{1,i})}{q_{\alpha}(\mathbf{t}_{2,j})} \right] \\
&= \mathbb{E} \log \left[1 + \frac{q_{\alpha}(\mathbf{t}_{2,i})}{p_{\beta}(\mathbf{t}_{2,i}|\mathbf{t}_{1,i})} (N-1) \right] \\
&\geq \mathbb{E} \log \left[\frac{q_{\alpha}(\mathbf{t}_{2,i})}{p_{\beta}(\mathbf{t}_{2,i}|\mathbf{t}_{1,i})} N \right]
\end{aligned}$$

We design our negative ‘‘proposal’’ distribution as $q_{\alpha}(\mathbf{t}_2) = \alpha_0 p(\mathbf{t}_2) + \sum_{k=1}^K \alpha_k p(\bar{\mathbf{v}}_2^k) p(\mathbf{v}_2^k)$. Insert this to the inequality above we obtain:

$$\begin{aligned}
& \mathcal{L}_{\text{TNCE}}^{\text{OPT}}(\alpha, \beta) \\
&\geq \mathbb{E} \log \left[\frac{\alpha_0 p(\mathbf{t}_{2,i}) + \sum_{k=1}^K \alpha_k p(\bar{\mathbf{v}}_{2,i}^k) p(\mathbf{v}_{2,i}^k)}{p_{\beta}(\mathbf{t}_{2,i}|\mathbf{t}_{1,i})} N \right] \\
&\geq \mathbb{E} \log \left[\frac{(p(\mathbf{t}_{2,i}))^{\alpha_0} \prod_{k=1}^K (p(\bar{\mathbf{v}}_{2,i}^k) p(\mathbf{v}_{2,i}^k))^{\alpha_k}}{p_{\beta}(\mathbf{t}_{2,i}|\mathbf{t}_{1,i})} N \right] \\
&= \mathbb{E} \log \left[\frac{p(\mathbf{t}_{2,i})}{p_{\beta}(\mathbf{t}_{2,i}|\mathbf{t}_{1,i})} \prod_{k=1}^K \left(\frac{p(\bar{\mathbf{v}}_{2,i}^k) p(\mathbf{v}_{2,i}^k)}{p(\mathbf{t}_{2,i})} \right)^{\alpha_k} N \right] \\
&\approx \mathbb{E} \log \left[\frac{p_{\beta}(\mathbf{t}_{2,i})}{p_{\beta}(\mathbf{t}_{2,i}|\mathbf{t}_{1,i})} \prod_{k=1}^K \left(\frac{p(\bar{\mathbf{v}}_{2,i}^k) p(\mathbf{v}_{2,i}^k)}{p(\mathbf{t}_{2,i})} \right)^{\alpha_k} N \right] \\
&= \log(N) - I(\mathbf{t}_{2,i}; \mathbf{t}_{1,i}|\beta) - \sum_{k=1}^K \alpha_k I(\mathbf{v}_{2,i}^k; \bar{\mathbf{v}}_{2,i}^k)
\end{aligned}$$

Therefore $I(\mathbf{t}_{2,i}; \mathbf{t}_{1,i}|\beta) + \sum_{k=1}^K \alpha_k I(\mathbf{v}_{2,i}^k; \bar{\mathbf{v}}_{2,i}^k) \geq \log(N) - \mathcal{L}_{\text{TNCE}}^{\text{OPT}}(\alpha, \beta)$.

F. Examples for negative sample optimization

This paper sets out with the aim of assessing the importance of complementary synergies between modalities. What is surprising is that increasing the ratio of weak modalities which is considered to contain less useful information can produce better-fused multi-modal representations. To be specific, the negative sample sampling ratio is roughly 1(RGB): 2(depth): 3(normal) in NYUv2 semantic segmentation task, 1(text): 1(video): 4(audio) in MOSI/MOSEI sentiment analysis task, and 1(point cloud): 3(RGB color): 2(height) in SUNRGB-D 3D Object Detection task. These results further support our idea of contrasting multi-modal

input tuples to avoid weak modalities being ignored and to facilitate modality fusion. A possible explanation for this might be that our proposed TupleInfoNCE encourages neural networks to use complement information in weak modalities, which can avoid networks converging to a lazy suboptimum where the network relies only on the strongest modality. This observation may support the hypothesis that mining useful information contained in weak modalities instead of always emphasizing strong modal information can obtain better fusion representations.

G. Examples for positive sample optimization

As mentioned in the main paper, we need to optimize the positive samples so that the contrastive learning can keep the shared information between positive and anchor samples while abstracting away nuisance factors. In the NYUv2 semantic segmentation task, rotation of 40 degrees, cropping with 160 center pixels, and Gaussian noise with a variance of 50 is best for RGB modality, while no augmentation is better for depth and normal modalities. As for the SUNRGB-D 3D Object Detection task, the final augmentation is to rotate the point cloud by 10 degrees, add Gaussian noise with a variance of 30 to the RGB image, and apply Gaussian noise with a variance of 10 to the height modality.

H. Dropout training

Dropout training Note $\mathbf{g}^{\alpha\beta}$ is originally designed to consume all input modalities, while in the crossmodal discrimination task, $\mathbf{g}^{\alpha\beta}$ needs to handle missing modalities as shown in Equation. We adopt a simple *dropout training* strategy to achieve this goal. To be specific, we randomly mask out modalities and fill them with placeholder values in the input. The missing modalities are the same in the positive and negative samples, yet could be different in the anchor tuple. This dropout strategy is only adopted with a probability of 0.6 for each training batch and the rest of the time we feed complete inputs to the feature encoder.

Robustness to dropout training We found dropout training does not cause a drop in feature quality compared to non-dropout training. We first obtain the optimal hyper-parameter by grid searching. When we use a dropout training strategy with fixed optimal hyper-parameter in NYUv2 semantic segmentation task, we achieve 48.3 mIoU performance which only outperforms our strategy by 0.2%, which presumably means dropout training might distill information from other modalities to avoid representation quality degradation.

I. Implementation Details

For NYU-Depth-V2 semantic segmentation task, we train the model with one V100 GPU for 500 epochs. The batch size is 20. We use SGD+momentum optimizer with an initial

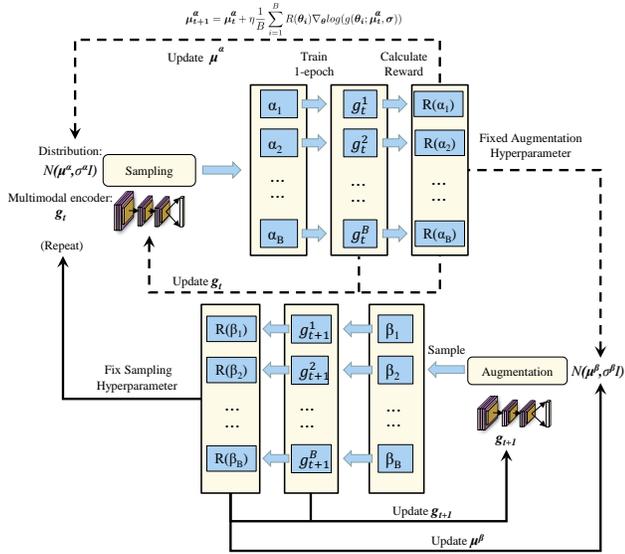


Figure 7. Overview of sample optimization.

learning rate 0.25. We use onecycle learning rate scheduler. We set the weight decay as $1e-4$. In the experiment of SUN RGB-D 3D object detection, our settings are consistent with VoteNet. We train the model on one 2080Ti GPU for 180 epochs. The initial learning rate is 0.001. We sample 20,000 points from each scene and the voxel size is 5cm. As for sentiment analysis on MOSEI and MOSI, We used the same data as MISA, where text features are extracted from BERT. Batch size is 32 for MOSI, and 16 for MOSEI. For sample optimization details, we provide a flow chart Figure 7 to further support the main paper.