

# Joint Wasserstein Autoencoders for Aligning Multimodal Embeddings

Shweta Mahajan    Teresa Botschen    Iryna Gurevych    Stefan Roth

Department of Computer Science, TU Darmstadt

## Abstract

One of the key challenges in learning joint embeddings of multiple modalities, e.g. of images and text, is to ensure coherent cross-modal semantics that generalize across datasets. We propose to address this through joint Gaussian regularization of the latent representations. Building on Wasserstein autoencoders (WAEs) to encode the input in each domain, we enforce the latent embeddings to be similar to a Gaussian prior that is shared across the two domains, ensuring compatible continuity of the encoded semantic representations of images and texts. Semantic alignment is achieved through supervision from matching image-text pairs. To show the benefits of our semi-supervised representation, we apply it to cross-modal retrieval and phrase localization. We not only achieve state-of-the-art accuracy, but significantly better generalization across datasets, owing to the semantic continuity of the latent space.

## 1. Introduction

The availability of significant amounts of image-text data on the internet (e.g., images with their captions) has posed the question whether it is possible to leverage information from both visual *and* textual sources. To take advantage of such heterogeneous data, one of the fundamental challenges is the joint representation of multiple domains [3]. Powerful multimodal representations are integral to the accuracy of models on cross-domain tasks, such as image captioning [21] or cross-domain retrieval [1, 7, 9, 24, 26, 43].

Multimodal embeddings of images and texts can be obtained by mapping input image and text representations into a common latent space [9, 43]. Learning such representations is often formulated as an image-text matching problem in a fully supervised setup. An alternative approach is to first learn separate latent spaces and to align them later through constraints, e.g., supervised information [40]. The benefit is that the latent representations of each modality are learned independently, allowing to take advantage of unsupervised (i.e. unpaired) data.

One of the main challenges in multimodal learning is to obtain meaningful latent representations such that they cap-

ture semantics that are present across modalities, even if the paired training data does not extensively cover all relevant semantic concepts. For example, the Flickr30k [48] and COCO [27] image captioning datasets do not contain matching image pairs. Lacking within-domain structural constraints, semantic similarity within each modality may not be preserved in the embedding space (Figs. 1b and 2a).

We address this problem by learning *semantically continuous* latent representations of images and texts in their respective embedding spaces, i.e. multimodal embeddings that encourage a smooth change in the semantics of the input modalities. We adopt a semi-supervised setting and propose a *joint Wasserstein autoencoder* (jWAE) model, leveraging that regularized autoencoders are known to yield semantically meaningful latent spaces [22, 39]. Specifically, we adopt Gaussian regularization to ensure semantically continuous latent representations of each input modality (c.f. Fig. 2b). In contrast to standard Wasserstein autoencoders (WAEs) [39], we share the Gaussian prior across modalities to encourage comparable levels of semantic continuity in both modalities. Unlike variational autoencoders (VAEs) [22], Wasserstein autoencoders map the input data to a point in the latent space, which allows for the coordination of the two modalities through a supervised loss based on matching image-text pairs. The advantage of the shared Gaussian prior on the two modalities is that their semantic representations are comparable and can be better aligned with supervision as illustrated in Fig. 1c.

We first evaluate our multimodal embeddings of images and texts on cross-modal retrieval and show that they yield state-of-the-art accuracy on the Flickr30k and COCO datasets. One of the crucial advantages of the semantically continuous representation from our semi-supervised approach is its generalization capability across datasets. The benefit over the state of the art widens when embeddings of one dataset are evaluated on a related, *previously unseen* dataset. Finally, we demonstrate the advantage of our jWAE on phrase localization on the Flickr30k Entities dataset [34], where we again outperform recent methods from the literature.

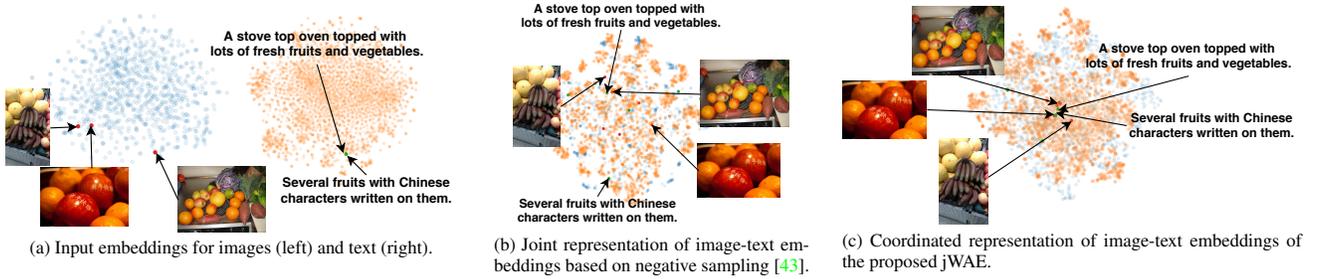


Figure 1. *From input embeddings to coordinated embeddings with semantic continuity:* (a) Semantic similarity in the input modalities does not imply proximity in the embedding space. (b) Joint representations align the modalities, but the sparsity of supervised information does not achieve continuity of the semantic space. (c) Our jWAE based on joint Gaussian regularization leads to semantic continuity.

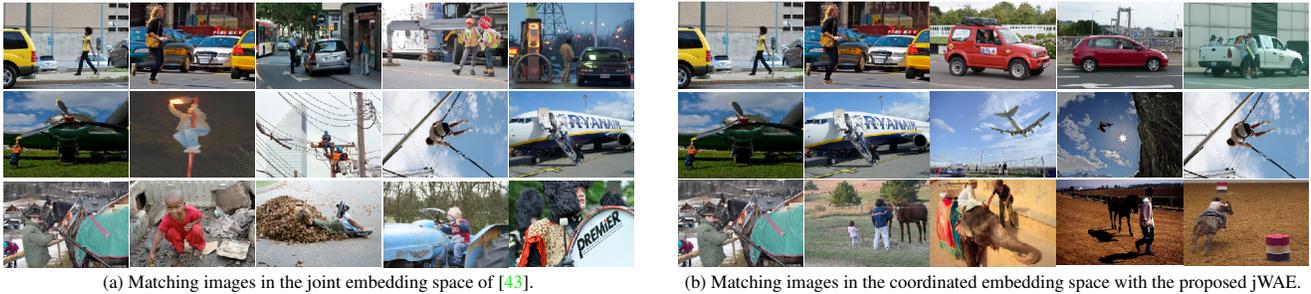


Figure 2. *Semantic continuity with Gaussian regularization.* Each row shows images that are close in the learned embedding space. (a) Without prior knowledge of matching image pairs, supervised cross-modal methods do not achieve semantic continuity. (b) Our jWAE based on Gaussian regularization achieves semantic continuity in the coordinated image space in absence of within-domain supervision.

## 2. Related Work

**Supervised multimodal learning.** Early work on multimodal learning includes Canonical Correlation Analysis (CCA) [17] and kernel CCA (KCCA) [25], which maximize correlation to learn projections of joint embeddings (*e.g.*, of images and texts). However, these methods do not scale to large datasets [28]. Deep Canonical Correlation Analysis (DCCA) [1] aims to overcome this scalability issue. Yet, optimization is challenging as the covariance matrix has to be estimated during training and is prone to over-fitting.

Many recent works formulate embedding multiple domains into a joint space as a learning-to-rank problem [9, 10, 14, 15, 23, 43, 46]. A ranking hinge loss with a margin is used such that matching cross-domain pairs are ranked higher (*i.e.* are closer in the latent space) than non-matching pairs. Wang *et al.* [43] additionally incorporate structural information on the input representations themselves with a domain-specific ranking loss. This requires prior knowledge about within-domain matching pairs, which was only available for text. Image-text matching has also been studied in a classification setting, employing logistic regression or a softmax with cross-entropy [11, 43].

Gu *et al.* [14] augment the ranking loss with a conditional generative model framework for cross-modal generation to obtain fine-grained multimodal features. Harada *et al.* [16] learn an image-text embedding space with a

Gaussian prior using generative adversarial networks. However, the Gaussian latent prior is applied only on the image modality, and the text distribution is matched to the latent image distribution. Moreover, the adversarial framework can suffer from mode collapse. Chi *et al.* [6] match image and text embeddings to the label representation. They assume that image and text have same label, which is limited to tasks where only one concept is required per image.

Wehrmann *et al.* [45] improve sentence representations with a character level inception module and [20, 26] improve image representations for image-text matching models. Huang *et al.* [20] use multi-label classification to extract various concepts in images, requiring additional image annotations. Lee *et al.* [26] propose an attention mechanism for aligning image regions with words in a sentence. This is orthogonal to the underlying multimodal embedding and can be combined with the proposed jWAE framework.

**Semi- and unsupervised multimodal learning.** Various approaches [31, 40] use autoencoders to obtain latent representations. Embeddings of the two domains can be aligned using distribution matching constraints [40, 42]. Unlike previous work, which does not encourage any continuity in the latent semantic space across modalities, we use *regularized* autoencoders based on generative models, which enforce the latent embeddings from the encoders to match a prior distribution, thereby yielding continuity in the embed-

Text	Supervised approach [43]	Ours (jWAE-MSE)
Two tan dogs play on the grass near the wall.		
A tennis player wearing white is jumping up to hit a ball.		

Table 1. Examples of images retrieved for a given sentence on Flickr30k based on embeddings trained on COCO.

ding space. We use the latent representations obtained from these models as a basis for our multimodal approach.

**Deep generative models** aim to minimize the difference between the model and the empirical data distribution, and have been successfully applied, *e.g.*, to image generation tasks. Generative adversarial networks (GANs) [13] generate the model distribution in a one step procedure where decoders are input with random (Gaussian) noise to construct the data distribution. Regularized autoencoders such as variational autoencoders (VAEs) [22] model the data distribution through a two step procedure. The empirical distribution is first mapped to a latent space via encoders and then mapped back to the data space via decoders. VAEs minimize the reconstruction error between the input and output representation and balance this with the discrepancy between the encoded representation in the latent space and a prior distribution, *e.g.* a Gaussian, for *each* input. Recently, [29, 39] proposed autoencoder-based frameworks where the discrepancy between the encoded distribution of *all* input representations and prior distribution is minimized. This forces the entire encoded distribution to match the prior. Such regularization captures the semantics of entire input distribution in a continuous latent representation, desirable for encoding each modality in multimodal learning.

### 3. Motivation & Background

Current approaches formulate the task of learning multimodal representations of images and text in an image-text matching framework [9, 44], in which a ranking loss is minimized in a fully supervised setting. For matching image-text pairs  $(x_l, y_l)$  and non-matching images  $x_{l'}$  or texts  $y_{l'}$  with similarity function  $s(x, y)$ , the ranking formulation based on the max-margin hinge loss with margin  $m$  is de-

Image	Supervised approach [43]	Ours (jWAE-MSE)
	<ul style="list-style-type: none"> <li>• Chevrolet car on display at a convention.</li> <li>• Construction in the city at night.</li> <li>• Two firemen beside their fire engine.</li> <li>• An escalator with many people on it, leading out of a tunnel.</li> </ul>	<ul style="list-style-type: none"> <li>• Two firemen beside their fire engine.</li> <li>• BSR trucks and machinery and workers.</li> <li>• An old, beat-up jeep being towed away.</li> <li>• Construction in the city at night.</li> </ul>
	<ul style="list-style-type: none"> <li>• The dog is running around the cow.</li> <li>• A person laying on the ground next to a cow.</li> <li>• Two children pet horses in a field.</li> <li>• Girl atop horse that is chasing a small longhorn.</li> </ul>	<ul style="list-style-type: none"> <li>• A lioness is chasing a black bison across a grassy plain.</li> <li>• A lioness chases a black animal with horns.</li> <li>• Two brown dogs play.</li> <li>• Horse jockeys racing on horses in a race.</li> </ul>

Table 2. Examples of the top-4 captions retrieved for a given image on Flickr30k based on embeddings trained on COCO.

defined as

$$\mathcal{L}_{MH} = \sum_l \psi \left( \max [0, m + s(x_l, y_{l'}) - s(x_l, y_l)] \right) + \psi \left( \max [0, m + s(x_{l'}, y_l) - s(x_{l'}, y_{l'})] \right). \quad (1)$$

Here,  $\psi$  is either the sum of hinge losses over all the negative samples for a given matching pair, or the maximum over all hinge losses. The dependence on the choice of negative examples limits the robustness and generalization of the obtained multi-modal embeddings; they fit to the particularities of a dataset, which can be seen from the example retrievals in Tables 1 and 2 on the Flickr30k dataset for an embedding space trained on the COCO dataset. While methods like domain adaptation rely on data from source *and* target datasets to adapt a model to perform better on a specific target dataset [19, 41], in this work we show the performance of the embeddings where both supervised and unsupervised losses are trained *only* on the source dataset, commonly referred to as *generalization*.

To overcome the limitation of existing methods in expressing semantic coherence within a domain and across multiple domains, we propose to employ Gaussian regularization on the latent distribution. By virtue of the autoencoder framework, structurally similar input representations are close to each other in the low-dimensional latent space; Gaussian regularization encourages *continuity* in the space of encoded representations. Semantic alignment of these spaces is further obtained with supervision. We refer to the resultant embeddings as *semantically continuous*. Our model consists of three main components. First, each input distribution is mapped to a Gaussian distribution where semantically similar representations are close to each other within the domain. This is illustrated in Fig. 2b, where images that are close in the embedding space are semantically

related. Second, we share this Gaussian prior across domains, leading to compatible levels of continuity in both domains. Third, the latent representations of images and text are semantically aligned with supervised information. As shown in Fig. 1c, images and texts representing similar semantics come closer in the proposed joint embedding space. Moreover, this offers significantly better generalization across datasets as seen in Tables 1 and 2, where the example captions retrieved with our semi-supervised approach are semantically more related to the given image than when only employing supervised ranking.

## 4. Approach

We propose a semi-supervised approach for improving the semantic alignment between two modalities, such as images and texts. Let  $X = \{x_i\}_{i=1}^{N_X}$  and  $Y = \{y_j\}_{j=1}^{N_Y}$  denote unpaired input images and texts, respectively. Further, let  $S = \{(x_l, y_l)\}_{l=1}^{N_S}$  be matching image and text pairs. We assume the latent space of each domain to be of dimension  $d$ . Encoders  $f_v : \mathbb{R}^{d_v} \rightarrow \mathbb{R}^d$  and  $f_t : \mathbb{R}^{d_t} \rightarrow \mathbb{R}^d$  map visual data (images) and text to their respective  $d$ -dimensional latent spaces.  $g_v : \mathbb{R}^d \rightarrow \mathbb{R}^{d_v}$  and  $g_t : \mathbb{R}^d \rightarrow \mathbb{R}^{d_t}$  are the decoders that map the latent representations back to images and text. We denote the latent representations of images and texts by  $\tilde{v}$  and  $\tilde{t}$ , respectively. Next, we let  $P_L \sim \mathcal{N}(0, I_d)$  be a unit Gaussian prior in the  $d$ -dimensional space with identity covariance matrix  $I_d$  and denote the encoded image and text distributions as  $F_v = \{f_v(x_i)\}_{i=1}^{N_X}$  and  $F_t = \{f_t(y_j)\}_{j=1}^{N_Y}$ , respectively. Similarly, we define the output (reconstructed) model distributions as  $G_v = \{g_v(\tilde{v}_i)\}_{i=1}^{N_X}$  and  $G_t = \{g_t(\tilde{t}_j)\}_{j=1}^{N_Y}$ . Following the common abuse of notation, we let  $F_v$  and  $F_t$  denote both the encoded activations and the distribution over encodings.

### 4.1. Wasserstein autoencoder backbone

We build the latent representation of each domain, images or text, on Wasserstein autoencoders [39]. We first describe the WAE backbone for the image pipeline and then extend it to text.

Generative models such as VAEs and WAEs minimize the discrepancy between the true data distribution  $X$  and the model (reconstructed) distribution  $G_v$ . In such models, latent variables  $z$ , sampled from a fixed prior distribution  $P_L$  in the latent space, are mapped to the space of the original data with parameterized functions  $g_v$ , and  $f$ -divergences, such as the KL-divergence or the Jensen-Shannon divergence, between the distributions are minimized. In high-dimensional spaces, estimating the model distribution directly by sampling in the latent domain would require a large number of samples from the latent distribution. Therefore, representing the model distribution by random sampling of the latent distribution is computationally expensive.

Variational autoencoders make sampling efficient by introducing a proposal distribution for each  $x_i$ . Specifically,  $F_v(x_i)$  is a latent distribution that generates latent representations  $z$  likely to produce  $x_i$ . The VAE minimizes

$$\mathcal{L}_{\text{VAE}} = \mathcal{D}_{\text{KL}}(F_v(x_i)||P_L) - \mathbb{E}_{F_v(x_i)}[\log G_v(\tilde{v}_i)], \quad (2)$$

where the first term encourages the latent variables  $z$  over each  $x_i$  to match the prior distribution  $P_L$ . While this helps latent samples to be representative of a data point, it does not capture the full underlying true data distribution. As has been pointed out by [39] for a Gaussian prior, this results in overlapping Gaussians in the latent space from different input data points. This also is the cause of blurry images from VAEs in image generation tasks. Moreover, mapping the input to a latent *distribution* is problematic when using supervision in the latent space, as intended here.

To model the entire data distribution in the latent space, minimizing  $\mathcal{D}(F_v||P_L)$  with  $F_v = \int F_v(x) dX$  is thus desirable [39]. To that end, given the input distribution  $X$  and the model distribution  $G_v$ , Wasserstein autoencoders (WAEs) minimize the optimal transport cost  $W_c(X, G_v)$  between the two distributions. Optimal transport with a cost function  $c(s, t) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$  is defined as

$$W_c(X, G_v) = \inf_{\gamma \in \Gamma(X, G_v)} \mathbb{E}_{(s,t) \sim \gamma} [c(s, t)], \quad (3)$$

where  $\Gamma(X, G_v)$  is the set of all possible joint distributions (couplings) of  $(s, t)$  whose marginals are  $X$  and  $G_v$ , respectively. In generative models with a deterministic mapping from the latent distribution  $P_L$  to  $G_v$ , the optimal transport cost between  $X$  and  $G_v$  reduces to finding a conditional distribution  $F_v$  such that  $\int F_v(x) dX = P_L$  [4, 39]. This constraint is enforced as a regularization term by minimizing  $\mathcal{D}(F_v||P_L)$ . This yields Wasserstein autoencoders as

$$\mathcal{L}_{\text{WAE}} = \inf_{F_v \in \mathcal{F}} \mathbb{E}_X \mathbb{E}_{F_v} [c(x, g_v(x))] + \lambda \mathcal{D}(F_v||P_L). \quad (4)$$

Here,  $c(s, t) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$  is a cost function and  $\mathcal{F}$  is a set of probabilistic encoders without any constraints (non-parametric), which model  $p(z|x)$ . We take it to be the set of all fully-connected networks with fixed size.

Choosing  $\mathcal{D}$  as the KL-divergence between  $F_v$  and  $P_L$  would require closed forms of  $F_v$  and  $P_L$ . Since the closed form of  $F_v$  is not available, we instead minimize the Jensen-Shannon divergence between the prior and encoded latent distributions  $\mathcal{D}_{\text{JS}}(F_v||P_L)$  using a GAN-based formulation, which allows to conveniently use samples from the distribution. We refer to  $\mathcal{D}_{\text{JS}}$  as Gaussian regularization since we consider the prior to be a unit Gaussian.

Note that the popular Wasserstein GANs [2] minimize the optimal transport cost (Eq. 3) between the data distribution and the model distribution using the dual formulation of the optimal transport cost directly from the latent space

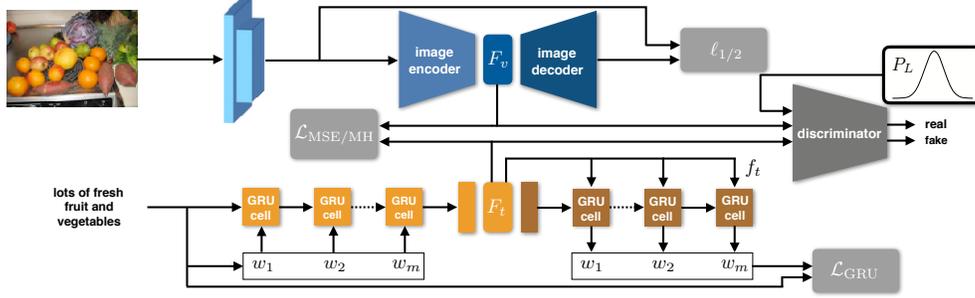


Figure 3. Architecture of our joint Wasserstein Autoencoder.

to the data space. That is, they do not have an encoder that maps the input to latent representations. However, for learning coordinated representations of two input data distributions with a shared prior, an encoder that maps data points to the latent space is desirable. We thus build on WAEs here.

## 4.2. Joint WAE

For learning coordinated representations, we are now interested in formulating continuous  $d$ -dimensional embedding spaces for each modality, which are aligned through constraints. To this end, we propose to share the prior on the latent representations. Specifically, the latent representations of each domain are constrained to be close to the same unit Gaussian distribution; that is we minimize the discrepancy  $\mathcal{D}(F_v||P_L)$  and  $\mathcal{D}(F_t||P_L)$  between the encoded representations of each modality and the Gaussian prior  $P_L$ .

We now formulate our joint Wasserstein autoencoder in terms of a classical encoder-decoder setting. To that end, we first note that when  $f_v$  and  $g_v$  in a regular WAE (Eq. 4) are parameterized with encoders and decoders in a deep neural network framework, the first term of Eq. (4) reduces to minimizing the reconstruction error between the true input representation and the decoded representation.

To apply the WAE formulation to sentences, we first extend it to gated recurrent units (GRUs) [37] where  $f_t$  is the encoded output of the GRU encoder and Gaussian regularization is applied to the encoded distribution,  $F_t$ , of all the sentences.  $f_t$  serves as the input hidden state to the GRU decoder and the reconstruction loss for the sentence is

$$\mathcal{L}_{GRU} = - \sum_{m=0}^{M-1} \log p_{g_t}(w_m^j | w_{0:m-1}^j, f_t(y_j); g_t), \quad (5)$$

where  $w_m^j$  is the ground truth word and  $p_{g_t}(w_m^j | w_{0:m-1}^j, f_t)$  is the output probability of word  $w_m^j$  in sentence  $y_j$  given the decoder  $g_t$  and hidden state  $f_t$ .

Recalling that  $F_v$  and  $F_t$  denote the encoded distribution of images and text, respectively, we formulate the unsuper-

vised part of our joint WAE loss as

$$\begin{aligned} \mathcal{L}_{jWAE} = & \lambda_1 \sum_{i=1}^{N_x} \|x_i - g_v(f_v(x_i))\| \\ & - \lambda_2 \sum_{j=1}^{N_y} \sum_{m=0}^{M-1} \log p_{g_t}(w_m^j | w_{0:m-1}^j, f_t(y_j); g_t) \\ & + \lambda_3 \mathcal{D}_{JS}(F_v || P_L) + \lambda_4 \mathcal{D}_{JS}(F_t || P_L). \end{aligned} \quad (6)$$

Here,  $\{\lambda_1, \dots, \lambda_4\}$  are the regularization parameters and  $\|\cdot\|$  is  $\ell_1$  or  $\ell_2$  norm. Modality specific reconstruction terms are crucial components in preventing mode collapse and encouraging diversity in the latent representations of each domain. For sentences encoded with fully connected encoders using pre-trained sentence encodings like average of word2vec [30], we use an analogous formulation with the mean-squared error as the reconstruction loss.

Gaussian regularization itself does not induce any semantic coupling between the cross-modal distributions. To ensure that the latent spaces not only have a compatible continuity but rather align semantic concepts across modalities, we add a supervised loss minimizing the distance between latent representations of matching image-text pairs. Embeddings of the two domains can now be directly aligned with the mean-squared error

$$\mathcal{L}_{MSE} = \frac{1}{N_S} \sum_{l=1}^{N_S} \|f_v(x_l) - f_t(y_l)\|^2. \quad (7)$$

The overall loss function of our approach is then given as

$$\mathcal{L}_{jWAE-MSE} = \mathcal{L}_{jWAE} + \mathcal{L}_{MSE}. \quad (8)$$

Alternatively, we also show the effect of Gaussian regularization with the max-margin hinge loss [9] from Eq. (1) as supervised loss function

$$\mathcal{L}_{jWAE-MH} = \mathcal{L}_{jWAE} + \mathcal{L}_{MH}. \quad (9)$$

The overall model architecture is illustrated in Fig. 3.

**Implementation.** The outputs from the encoders (*i.e.* the encoded distributions  $F_v$  and  $F_t$ ) along with  $z \sim P_L$  are

input to the discriminator. The number of samples from the prior distribution  $P_L$  equals the sum of the samples output by the two encoders. The discriminator distinguishes between the encoded distributions and the joint Gaussian prior. Note that we implement a single discriminator network for two generator networks, which makes our architecture computationally efficient. This is possible as the same prior distribution is used for images and texts. The discriminator is a fully-connected three layer neural network with leaky ReLU non-linearities after the first two layers, which enables a better flow of gradients during optimization [36]. The generator network (the encoder) of each pipeline tries to “fool” the discriminator network by generating encodings close to the Gaussian prior distribution.

In general, both encoders and decoders consist of two fully connected layers; ReLU non-linearities are applied after the first layers. For the text pipeline based on GRUs, the encoder is a bi-directional GRU with two layers. The output of the GRU is encoded in the latent space after application of a linear fully-connected layer. The decoder is a uni-directional GRU with word dropout encouraging meaningful latent representations of sentences.

## 5. Experiments

To show the applicability of our jWAE framework across *different tasks*, we evaluate the learned multimodal embeddings on cross-modal retrieval and phrase localization.

### 5.1. Cross-modal retrieval

**Visual input.** We consider pretrained VGG-19 [38] and ResNet-152 [18] models for image features. For VGG-19, we extract the 4096-dimensional feature vector from the first fully connected layer and for ResNet-152 the 2048-dimensional activations from the fully connected layer.

**Textual input.** Following [43], we use the mean of 300-dimensional word2vec [30] features of the words in the sentence. Alternatively, we use nonlinear Fisher vectors from a hybrid Gaussian-Laplacian mixture model (HGLMM) [24]. For GRU, one-hot encodings of the words are projected with an embedding layer, which is initialized either randomly or with pre-trained word2vec embeddings.

**Datasets.** Two popular benchmark datasets for evaluating multimodal visual and textual representations are Flickr30k and COCO [9, 24, 43]. Flickr30k [48] is comprised of 31783 images with five captions per image. We use the splits of [21, 44]; validation and test splits consist of 1000 images with 5 captions each. The remaining images are used for training. COCO [27] is larger and more diverse than Flickr30k. It consists of 82783 training images with five captions each. 5000 images from the validation set are retained for validation purposes and the remaining 30504

images are used for training. Similar to [21, 24, 44], we use 1000 images with their captions in the test split.

**Network training.** We train the network (see Appendix B for architectural details) using the Adam optimizer with a learning rate of 1E-4. The discriminator is trained with a learning rate of 5E-5. The batch size is taken as 64 or 128. For jWAE-MSE, the regularization parameters are set to  $\lambda_1 = \lambda_2 = 1.0$  for the reconstruction terms in Eq. (6) and  $\lambda_3 = \lambda_4 = 0.2$  for the Gaussian regularization. For jWAE-MH, the parameters are  $\lambda_1 = 0.5$ ,  $\lambda_2 = 0.005$ ,  $\lambda_3 = \lambda_4 = 0.01$  for most of the experiments.

**Evaluation metric.** In cross-modal retrieval tasks, Recall@ $K$  is a standard performance measure and defined as the fraction of instances for which their ground truth is in the top- $K$  based on a similarity score (cosine similarity).

**Baselines & methods.** We compare the accuracy of our approach against several state-of-the-art methods for image-to-text and text-to-image retrieval, particularly the Embedding Network of [43, 44] and VSE++ [9], which use different formulations of the ranking loss, and the attention-based Stacked Cross Attention Network (SCAN) [26]. For evaluation with SCAN, we integrate the jWAE-MH framework with the best-performing setting on the respective dataset. To compare our approach against methods without negative sampling, we also include the Similarity Network of [43] and the Canonical Correlation Analysis (CCA) approach of [24]. For our jWAE framework, we demonstrate the effect of Gaussian regularization with supervision through the mean-squared error (jWAE-MSE, Eq. 8) as well as a margin-based hinge loss (jWAE-MH, Eq. 9). We extend jWAE for the SCAN method (jWAE-MH+SCAN t-i/i-t) where in i-t attention is applied on words with respect to each image region and in t-i image regions are attended with respect to each word in the sentence. We also include results from a MMD loss for matching text and image distributions (MMD-MSE) [40] and an ablation of our method without the Gaussian regularization (MSE), both with the usual reconstruction error for autoencoders.

**Results.** In Table 3, we show the results of our method for cross-modal retrieval on the Flickr30k and COCO datasets.

Our jWAE framework leads to competitive results compared to the current state of the art in image-to-text retrieval. For example, jWAE-MH outperforms VSE++ with respect to top-1 recall by 1.0% points on Flickr30k and by 2.0% points on COCO. For the Embedding Network with VGG+w2v features, jWAE-MH has 3.5% better accuracy on COCO. Our semi-supervised representations also improve the top-1 recall of SCAN t-i by 2.2% and 2.1% points on Flickr30k and COCO, respectively. For text-to-image retrieval, improving the top-1 recall is more challenging. Yet, we also achieve an improvement in top-1 recall of 1.3% and

Method (Features)		Recall on Flickr30k						Recall on COCO					
		Image-to-text			Text-to-image			Image-to-text			Text-to-image		
		@1	@5	@10	@1	@5	@10	@1	@5	@10	@1	@5	@10
Without Negative Sampling	CCA (VGG+HGLMM) [24]	34.4	61.0	72.3	24.4	52.1	65.6	37.7	66.6	79.1	24.9	58.8	76.5
	CCA (Mean Vector) [24]	24.8	52.5	64.3	20.5	46.3	59.3	33.2	61.8	75.1	24.2	56.4	72.4
	Sim. Network (VGG+HGLMM) [43]	16.6	38.8	51.0	7.4	23.5	33.3	30.9	61.1	76.2	14.0	30.0	37.8
	MMD-MSE (VGG+w2v) [40]	35.3	60.8	72.8	16.5	40.2	52.9	42.9	73.8	84.4	19.6	50.2	66.6
With Negative Sampling	Emb. Network (VGG+w2v) [43]	35.7	62.9	74.4	25.1	53.9	66.1	40.7	74.2	85.3	33.5	68.7	83.2
	Emb. Network (VGG+HGLMM) [43]	40.3	68.9	79.9	29.7	60.1	72.1	50.1	79.7	89.2	39.6	75.2	86.9
	2WayNet (VGG+HGLMM) [7]	49.8	67.5	–	36.0	55.6	–	55.8	75.2	–	39.7	63.3	–
	VSE++ (Resnet+GRU) [9]	52.9	80.5	87.2	39.6	70.1	79.5	64.6	90.0	95.7	52.0	84.3	92.0
	GXN (ResNet+GRU) [15]	56.8	–	89.6	41.5	–	80.1	68.5	–	97.9	56.6	–	94.5
	SCO (Resnet+GRU) [20]	55.5	82.0	89.3	41.1	60.5	80.1	69.9	92.9	97.5	56.7	<b>87.5</b>	<b>94.8</b>
	SCAN t-i (ResNet+GRU) [26]	61.8	87.5	93.7	45.8	74.4	83.0	70.9	<b>94.5</b>	97.8	56.4	87.0	93.9
	SCAN i-t (ResNet+GRU) [26]	67.9	89.0	94.4	43.9	74.2	82.8	69.2	93.2	97.5	54.4	86.0	93.6
	MSE (VGG+w2v)	33.9	61.3	73.2	15.9	40.3	52.8	45.0	76.5	87.0	21.3	52.3	68.1
	jWAE-MSE (VGG+w2v)	35.7	61.6	73.6	17.3	41.8	55.3	43.2	75.1	85.7	21.5	53.5	69.3
	jWAE-MH (VGG+w2v)	35.4	62.5	74.4	24.1	50.4	62.6	44.2	77.4	87.7	31.3	66.0	81.3
	jWAE-MSE (VGG+HGLMM)	40.3	66.3	77.2	20.3	46.5	58.9	50.3	79.4	88.3	25.2	57.5	73.3
	jWAE-MH (ResNet+GRU)	53.9	82.2	87.4	40.7	72.4	81.9	66.6	91.4	96.6	53.1	84.5	92.0
	jWAE-MH+SCAN t-i (ResNet+GRU)	64.0	89.4	<b>95.2</b>	<b>47.1</b>	<b>75.9</b>	<b>84.0</b>	<b>72.0</b>	<b>94.5</b>	<b>98.3</b>	<b>57.1</b>	<b>87.5</b>	94.1
	jWAE-MH+SCAN i-t (ResNet+GRU)	<b>68.5</b>	<b>89.6</b>	94.2	44.2	74.2	83.2	69.8	93.3	98.1	54.6	86.1	93.2

Table 3. Cross-modal retrieval results (in %) on the Flickr30k dataset [48] as well as the COCO dataset [27] with 1000 test images.

Method (all ResNet+GRU)	Image-to-text		Text-to-image	
	R@1	R@5	R@1	R@5
VSE++ [9] (baseline)	64.6	90.0	52.0	84.3
VAE-MH	57.8	86.9	46.3	80.3
jWAE-MH	<b>66.6</b>	<b>91.4</b>	<b>53.1</b>	<b>84.5</b>

Table 4. Comparison of jWAE with VAE on the COCO dataset.

Flickr30k(Train) $\Rightarrow$ COCO(Test)				
Method	Image-to-text		Text-to-image	
	R@5	R@10	R@5	R@10
CCA (VGG+w2v) [24]	13.3	20.1	10.3	16.0
Embed. Network (VGG+w2v) [43]	37.6	49.5	32.5	45.8
MSE (VGG+w2v)	40.1	52.9	21.7	33.8
SCAN (ResNet+GRU) [26]	59.2	70.0	53.2	66.6
jWAE-MSE (VGG+w2v)	42.8	55.3	28.3	40.9
jWAE-MH (VGG+w2v)	43.4	58.4	34.5	49.2
jWAE-MH+SCAN (ResNet+GRU)	<b>64.4</b>	<b>76.9</b>	<b>55.2</b>	<b>68.7</b>

Table 5. Generalization results of models trained only on the Flickr30k training set and evaluated on the COCO test set.

0.7% points over SCAN t-i, and an improvement of 1.1% for VSE++ on Flickr30k and COCO. This shows that irrespective of the network architecture and complexity of input features, our jWAE improves the current state of the art. To the best of our knowledge, our jWAE framework improves (over) the currently leading cross-modal retrieval methods.

We observe that the recall for text-to-image retrieval of jWAE-MSE is not as competitive and is comparable to methods that do not use negative sampling. This can be attributed to the nature of the datasets where five sentences

compete for the same image. Moreover, given a sentence there can be many images that can be described reasonably well by the sentence. jWAE-MH bridges the gap between maximizing top-K recall by ranking matching image-text pairs higher than non-matching pairs *and* improving accuracy of the embeddings by encouraging semantic continuity.

In Table 4, we additionally compare jWAE-based Gaussian regularization to traditional VAEs using VSE++ as baseline. While jWAE-MH improves the accuracy of the VSE++ baseline, VAE-MH results in a decrease of top-K recall. The reason is that VAEs map *each* input to a Gaussian latent *distribution*. This hinders the application of a point-wise supervised loss in the latent space. In order to match the latent representations of two modalities with the MH loss, we instead require mapping to a *point* in the latent space. jWAE enforces global Gaussian priors on the latent distributions while mapping the input to a point in the embedding space. Therefore and unlike VAEs, jWAEs are suitable for semi-supervised learning of joint embeddings.

To show that our method learns meaningful representations with continuity in the latent space, we test the cross-dataset generalization capability of our method against various retrieval approaches: CCA [24], the Embedding Network [43], and SCAN [26]. For cross-dataset generalization, the model is trained on the training set of one dataset, *e.g.* COCO (Flickr30k), and tested on the test set of another dataset, *i.e.* Flickr30k (COCO). Note, we do not train to improve the accuracy on specific target dataset. We find that previous methods based on global representations [43] have low generalization performance with top-10 recall as low as 12.2% when testing on Flickr30k for a model trained

COCO(Train)  $\Rightarrow$  Flickr30k(Test)

Method	Image-to-text		Text-to-image	
	R@5	R@10	R@5	R@10
Embed. Network (VGG+w2v) [43]	33.7	45.5	8.4	12.2
MSE (VGG+w2v)	43.6	56.9	24.3	34.2
SCAN (ResNet+GRU) [26]	73.7	82.6	61.3	72.4
jWAE-MSE (VGG+w2v)	48.5	60.0	29.1	40.7
jWAE-MH (VGG+w2v)	51.1	63.3	37.0	49.1
jWAE-MH+SCAN (ResNet+GRU)	<b>80.0</b>	<b>87.0</b>	<b>66.7</b>	<b>75.9</b>

Table 6. Generalization results of model trained only on the COCO training set and evaluated on the Flickr30k test set.

on COCO. Fine-grained representations based on attention [26] generalize better compared to [43]. Following Tables 5 and 6, the jWAE-MH framework significantly improves the generalization across datasets further, owing to the semantic continuity from the Gaussian regularization. For image-to-text and text-to-image retrieval we improve the top-5 recall by 5.2% and 2.0% points, respectively, generalizing from Flickr30k to COCO and by 6.3% and 5.4% points, respectively, for generalizing from COCO to Flickr30k.

In Fig. 2, we show modality-specific semantic continuity, where, *e.g.*, in the third row, for an image with “a man and a horse”, matching images retrieved by our approach are of ‘person-animal interaction’ whereas for the supervised approach [43] matching images show the concept ‘person’. Similarly, for cross-modal semantic continuity in Table 1, [43] retrieves an image with a concept ‘two’ while jWAE is able to retrieve the image representative of the given sentence “Two tan dogs play on grass near the wall”. We provide additional qualitative examples in Appendix C.

In the Appendix A, we additionally study the accuracy under limited supervision.

## 5.2. Phrase localization

We next analyze the benefit of our jWAE framework for phrase localization on the Flickr30k Entities dataset [34]. Phrase localization associates (grounds) a phrase to a region in the image using bounding boxes [5, 35, 43, 47]. Following [43], we formulate phrase localization as a retrieval problem where given an image and a phrase from its associated sentence, the phrase is mapped to the regions in the image. Bounding box proposal regions are extracted with Edge Box [49]. Since we are mainly interested in evaluating the quality of our multimodal embeddings rather than the specific task, we compared to other embedding-based approaches [35, 43]. Additionally, we integrate our jWAE framework with Conditional Image Text Embedding (CITE) [32], which builds on top of the embeddings from the Similarity Network. We also include [33], which uses additional image and language constraints, and [47], which considers all possible bounding boxes based on image concepts like segmentation, word priors, and detection scores.

Methods	R@1	R@5	R@10
MCB [11]	48.7	–	–
GroundER [35]	47.8	–	–
Embedding Network [43]	51.0	70.4	75.5
Similarity Network [43]	51.0	70.3	75.0
SPC [33]	55.4	–	–
IGOP [47]	53.9	–	–
CITE [32]	59.2	–	–
jWAE-MSE	52.5	<b>75.0</b>	<b>81.5</b>
jWAE-MSE+CITE	<b>60.4</b>	–	–

Table 7. Phrase localization on the Flickr30k Entities dataset [34].

**Dataset and input.** The Flickr30k Entities dataset [34] augments the captions of images in Flickr30k with 244k mentions of distinct entities across sentences. The mentions are associated with 276k bounding boxes. Similar to [32, 35, 43], we extract 4096-dimensional visual features from Fast R-CNN [12], finetuned on the PASCAL VOC 2007–2012 datasets [8]. We use proposal regions with IoU  $\geq 0.7$  as a positive region for a phrase during training. For encoding phrases, PCA is applied to HGLMM features to reduce the dimensionality to 6000 [24].

**Results.** We compare our method with [11, 35, 43] where multimodal embeddings are evaluated for phrase localization. Following these methods, we use 200 or 500 Edge Box proposals per image. An IoU of at least 0.5 is required for a proposal region to match the ground truth bounding box for a phrase. As shown in Table 7, our method outperforms previous multimodal embedding networks for the phrase localization task by 1.5% for top-1 recall. The gap compared to [43] widens to around 5% for the top-5 and top-10 recall. Moreover, using jWAE as the embedding framework in CITE [32] similarly improves the top-1 recall by 1.2% with new state of the art results for phrase localization. This again highlights the improved accuracy of the embeddings obtained from our semi-supervised jWAE approach. Please see Appendix D for additional qualitative results.

## 6. Conclusion

We presented a novel joint Wasserstein autoencoder framework for modeling continuous multimodal representations of images and texts with Gaussian regularization, allowing to better capture the semantic structure in latent representations. Our experiments show that our multimodal embeddings push the current state of the art under full supervision. A key advantage of our method is its generalization capability across datasets, where it significantly outperforms recent methods. We thus believe our semi-supervised approach provides an important step toward learning generalizable multimodal representations, which are a crucial component, *e.g.*, for image captioning [10] in the real world.

**Acknowledgement.** This work has been supported by the German Research Foundation as part of the Research Training Group *Adaptive Preparation of Information from Heterogeneous Sources (AIPHES)* under grant No. GRK 1994/1.

## References

- [1] Galen Andrew, Raman Arora, Jeff A. Bilmes, and Karen Livescu. Deep canonical correlation analysis. In *ICML*, 2013. 1, 2
- [2] Martín Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein GAN. *ICML*, 2017. 4
- [3] Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. *TPAMI*, 41(2), 2019. 1
- [4] Olivier Bousquet, Sylvain Gelly, Ilya Tolstikhin, Carl-Johann Simon-Gabriel, and Bernhard Schölkopf. From optimal transport to generative modeling: The VEGAN cookbook. *arXiv:1705.07642*, 2017. 4
- [5] Kan Chen, Rama Kovvuri, Jiyang Gao, and Ram Nevatia. MSRC: Multimodal spatial regression with semantic context for phrase grounding. In *ICMR*, 2017. 8
- [6] Jingze Chi and Yuxin Peng. Dual adversarial networks for zero-shot cross-media retrieval. In *IJCAI*, 2018. 2
- [7] Aviv Eisenschat and Lior Wolf. Linking image and text with 2-way nets. In *CVPR*, 2017. 1, 7
- [8] Mark Everingham, Luc J. Van Gool, Christopher K. I. Williams, John M. Winn, and Andrew Zisserman. The Pascal visual object classes (VOC) challenge. *IJCV*, 88(2), 2010. 8
- [9] Fartash Faghri, David J. Fleet, Ryan Kiros, and Sanja Fidler. VSE++: Improving visual-semantic embeddings with hard negatives. In *BMVC*, 2018. 1, 2, 3, 5, 6, 7, 11
- [10] Andrea Frome, Gregory S. Corrado, Jonathon Shlens, Samy Bengio, Jeffrey Dean, Marc’Aurelio Ranzato, and Tomas Mikolov. Devise: A deep visual-semantic embedding model. In *NIPS*, 2013. 2, 8
- [11] Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. In *EMNLP*, 2016. 2, 8
- [12] Ross B. Girshick. Fast R-CNN. In *ICCV*, 2015. 8
- [13] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial networks. In *NIPS*, 2014. 3
- [14] Jiuxiang Gu, Jianfei Cai, Shafiq R. Joty, Li Niu, and Gang Wang. Look, imagine and match: Improving textual-visual cross-modal retrieval with generative models. In *CVPR*, 2018. 2
- [15] Guibing Guo, Songlin Zhai, Fajie Yuan, Yuan Liu, and Xingwei Wang. VSE-ens: Visual-semantic embeddings with efficient negative sampling. In *AAAI*, 2018. 2, 7
- [16] Tatsuya Harada, Kuniaki Saito, Yusuke Mukuta, and Yoshitaka Ushiku. Deep modality invariant adversarial network for shared representation learning. In *ICCV Workshops*, 2017. 2
- [17] David R. Hardoon, Sándor Szedmák, and John Shawe-Taylor. Canonical correlation analysis: An overview with application to learning methods. *Neural Computation*, 16(12), 2004. 2
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 6
- [19] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei A. Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *ICML*, 2018. 3
- [20] Yan Huang, Qi Wu, Chunfeng Song, and Liang Wang. Learning semantic concepts and order for image and sentence matching. In *CVPR*, 2018. 2, 7
- [21] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. *TPAMI*, 39(4), 2017. 1, 6
- [22] Diederik P. Kingma and Max Welling. Auto-encoding variational Bayes. In *ICLR*, 2014. 1, 3
- [23] Ryan Kiros, Ruslan Salakhutdinov, and Richard S. Zemel. Unifying visual-semantic embeddings with multimodal neural language models. *TACL*, 2015. 2
- [24] Benjamin Klein, Guy Lev, Gil Sadeh, and Lior Wolf. Associating neural word embeddings with deep image representations using Fisher vectors. In *CVPR*, 2015. 1, 6, 7, 8
- [25] Pei Ling Lai and Colin Fyfe. Kernel and nonlinear canonical correlation analysis. *Int. J. Neural Syst.*, 10(5), 2000. 2
- [26] Kuang-Huei Lee, Xi Chen, Gang Hua, Houdong Hu, and Xiaodong He. Stacked cross attention for image-text matching. In *ECCV*, 2018. 1, 2, 6, 7, 8, 11
- [27] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014. 1, 6, 7
- [28] Zhuang Ma, Yichao Lu, and Dean P. Foster. Finding linear structure in large datasets with scalable canonical correlation analysis. In *ICML*, 2015. 2
- [29] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, and Ian J. Goodfellow. Adversarial autoencoders. *ICLR Workshops*, 2016. 3
- [30] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *ICLR Workshops*, 2013. 5, 6
- [31] Jiquan Ngiam, Aditya Khosla, Mingyu Kim, Juhan Nam, Honglak Lee, and Andrew Y. Ng. Multimodal deep learning. In *ICML*, 2011. 2
- [32] Bryan A. Plummer, Paige Kordas, M. Hadi Kiapour, Shuai Zheng, Robinson Piramuthu, and Svetlana Lazebnik. Conditional image-text embedding networks. In *ECCV*, 2018. 8, 12, 13
- [33] Bryan A. Plummer, Arun Mallya, Christopher M. Cervantes, Julia Hockenmaier, and Svetlana Lazebnik. Phrase localization and visual relationship detection with comprehensive image-language cues. In *ICCV*, 2017. 8
- [34] Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and Svetlana Lazebnik.

- Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. *IJCV*, 123(1), 2017. [1](#), [8](#)
- [35] Anna Rohrbach, Marcus Rohrbach, Ronghang Hu, Trevor Darrell, and Bernt Schiele. Grounding of textual phrases in images by reconstruction. In *ECCV*, 2016. [8](#)
- [36] Tim Salimans, Ian J. Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training GANs. In *NIPS*, 2016. [6](#)
- [37] Mike Schuster and Kuldip K. Paliwal. Bidirectional recurrent neural networks. *TSP*, 45(11), 1997. [5](#)
- [38] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *ICLR*, 2015. [6](#)
- [39] I. Tolstikhin, O. Bousquet, S. Gelly, and B. Schölkopf. Wasserstein auto-encoders. In *ICLR*, 2018. [1](#), [3](#), [4](#)
- [40] Yao-Hung Hubert Tsai, Liang-Kang Huang, and Ruslan Salakhutdinov. Learning robust visual-semantic embeddings. In *ICCV*, 2017. [1](#), [2](#), [6](#), [7](#)
- [41] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *CVPR*, 2017. [3](#)
- [42] Bokun Wang, Yang Yang, Xing Xu, Alan Hanjalic, and Heng Tao Shen. Adversarial cross-modal retrieval. In *ACM MM*, 2017. [2](#)
- [43] Liwei Wang, Yin Li, Jing Huang, and Svetlana Lazebnik. Learning two-branch neural networks for image-text matching tasks. *TPAMI*, 41(2), 2019. [1](#), [2](#), [3](#), [6](#), [7](#), [8](#), [14](#)
- [44] Liwei Wang, Yin Li, and Svetlana Lazebnik. Learning deep structure-preserving image-text embeddings. In *CVPR*, 2016. [3](#), [6](#)
- [45] Jonas Wehrmann and Rodrigo C. Barros. Bidirectional retrieval made simple. In *CVPR*, 2018. [2](#)
- [46] Jason Weston, Samy Bengio, and Nicolas Usunier. WSA-BIE: Scaling up to large vocabulary image annotation. In *IJCAI*, 2011. [2](#)
- [47] Raymond Yeh, Jinjun Xiong, Wen-Mei W. Hwu, Minh Do, and Alexander G. Schwing. Interpretable and globally optimal prediction for textual grounding using image concepts. In *NIPS*, 2017. [8](#)
- [48] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *ACL*, 2014. [1](#), [6](#), [7](#)
- [49] C Lawrence Zitnick and Piotr Dollár. Edge boxes: Locating object proposals from edges. In *ECCV*, 2014. [8](#)

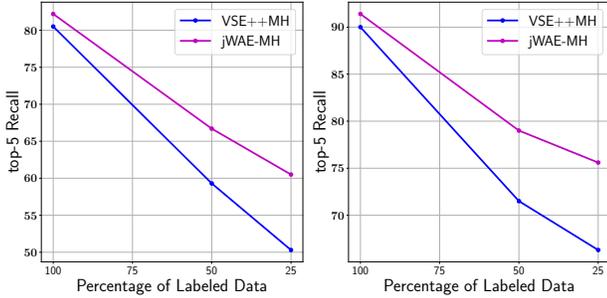


Figure 4. Comparison of a standard supervised loss (MH) and the proposed semi-supervised jWAE under limited supervision: (left) Flickr30k; (right) COCO.

## Appendix A. jWAE with Limited Supervision

To show the effectiveness of our jWAE in minimally supervised settings, we perform additional experiments with limited labeled data in the training set, specifically with 25% and 50% of the labeled data. If applicable, the rest of the data is included in an unpaired fashion for semi-/unsupervised learning. We observe that compared to standard supervised methods for learning joint embeddings of images and text [9], our semi-supervised approach performs better under different levels of supervision, highlighting the benefit of the semantic continuity from the joint Gaussian regularization. We observe that the accuracy gap widens as we decrease the supervision level, *e.g.* from 7.4% to 10.2% points as we decrease the supervision from 50% to 25% on Flickr30k (*c.f.* Fig. 4).

Table 8 shows the effect of training with limited supervision using region-based features and attention. Our variant of SCAN [26], *i.e.* SCAN+jWAE-MH, outperforms SCAN by 3.4% and 1.3% on COCO and Flickr30k, respectively. Since the basic SCAN is trained with 36 image regions per image, which effectively increases the number of training points, SCAN performs reasonably well even at 25% supervision. Adding our jWAE on top increases the robustness to limited supervision further. Moreover, further decreasing the supervision to 5%, the top-1 recall of SCAN drops to 42.0% while for SCAN+jWAE-MH the top-1 recall remains high at 54.3% on the COCO dataset. This highlights the improved generalization performance of our jWAE framework at different supervision levels.

We also compare against other methods that rely on additional unsupervised information, particularly the unsupervised autoencoder reconstruction loss combined with supervised MSE-based alignment, but without the Gaussian regularization (MSE). The top-5 recall of our jWAE-MSE method is higher compared to the semi-supervised MSE baseline by 2.6% and 3.2% points on average for Flickr30k at 50% and 25% supervision, respectively. Similarly for the COCO dataset, the top-5 recall is 2.8% and 3.1% points

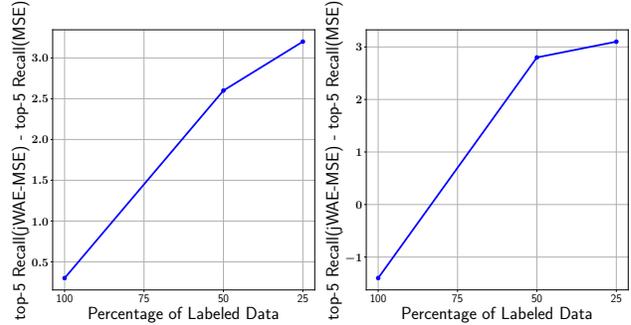


Figure 5. Percentage gain in top-5 recall with the proposed jWAE-MSE over the MSE baseline under limited supervision with 50% and 25% of the labeled data: (left) Flickr30k; (right) COCO.

Method(ResNet+GRU) (label %)	COCO		Flickr30k	
	R@1	R@5	R@1	R@5
SCAN [26] (25%)	55.9	86.5	45.9	72.6
SCAN [26] (5%)	42.0	84.3	–	–
SCAN+jWAE-MH (25%)	59.3	88.3	47.2	75.0
SCAN+jWAE-MH (5%)	54.3	87.0	–	–

Table 8. Comparison of SCAN and SCAN-jWAE under limited supervision.

higher compared to MSE at 50% and 25% supervision, respectively, as shown in Fig. 5. The difference in recall at a particular supervision level between our method and the baseline increases with decreasing labeled data. This shows that encouraging continuity in the latent spaces helps to better align similar encoded representations and can be harnessed in learning cross-modal concepts even when supervised information is limited.

## Appendix B. Network Architecture

### Cross-modal retrieval

**Image pipeline.** The input representations obtained from VGG-19 or ResNet-152 are fed into our joint Wasserstein autoencoder. The image encoder takes 4096 inputs (2048 for ResNet-152), which are fully connected to a hidden layer of 2048 nodes. The encoder outputs into a  $d$ -dimensional latent space. The decoder is symmetric to the encoder. We apply ReLU nonlinearities after the hidden layer in both encoder and decoder.

**Text pipeline.** We use a similar network with one hidden layer for the text encoder. For 300-dimensional sentence inputs, the hidden encoder layer is of dimensionality 1024, and the latent layer has 256 dimensions. The same latent dimensionality is chosen for the image pipeline. The decoder network is again symmetric to the encoder. We apply ReLU nonlinearities after the hidden layer in the encoder and the decoder. When using high-dimensional HGLMM features,

Text	Ground truth image	Retrieved image with jWAE-MSE
A man using his laptop computer while a cat sits on his lap.		
A pizza and grapes sit on a tray next to a drink.		
A bath room with a toilet a sink and a mirror.		
A kitchen with a stove, oven, and refrigerator.		
A slice of chocolate cake with dark chocolate icing.		
a no skate boarders sign on the side of the road.		

Table 9. Given a caption, examples of images retrieved by our method that did not match the ground truth. While they are clearly incorrect, they bear a semantic relation to the ground truth.

we set the hidden layer size to 3000 and use a 512 (1024)-dimensional latent space for jWAE-MSE (jWAE-MH).

For GRU, the dimension of word embeddings is 300. Starting with one-hot coded vectors, the GRU encoder consists of an embedding layer, a bidirectional GRU layer, and a fully connected layer with dimensionality 1800. The GRU decoder has a fully connected layer and a uni-directional GRU layer. Word dropout regulates the dependence of the prediction of a word in a sentence on the preceding ground truth sequence. To obtain a latent encoding that is representative of the sentence, we set the word dropout rate to 0.2 in the decoder.

**Discriminator network.** We use a three layer discriminator network on the latent spaces of each pipeline. The sizes of the layers are chosen as 256, 256, and 2, respectively. We use leaky ReLU activations between all the layers.

Image	Retrieved text with jWAE-MSE
	<ul style="list-style-type: none"> <li>• A couple of men that are walking around on some grass.</li> <li>• Two men are playing with a toy in a wooded area.</li> <li>• A couple of men standing on top of a lush green forest.</li> <li>• <b>A group of young men and women sitting at a table.</b></li> </ul>
	<ul style="list-style-type: none"> <li>• Painting of oranges, a bowl, candle, and a pitcher.</li> <li>• <b>Four boats with people carrying lots of bananas and other foods.</b></li> <li>• <b>A couple of people in boats with food.</b></li> <li>• Painting of a table with fruit on top of it.</li> </ul>
	<ul style="list-style-type: none"> <li>• A man riding skis down a snow covered slope.</li> <li>• A person on skis going down a snow covered hill.</li> <li>• A person turning while skiing down a snowy hill.</li> <li>• A lady snow skiing on flat ground</li> </ul>
	<ul style="list-style-type: none"> <li>• Some purple bananas and other fruits are together.</li> <li>• Some purple bananas sitting between apples and some squash.</li> <li>• A pile of black bananas and other fruit.</li> <li>• Assorted fruit on display at a fruit market.</li> </ul>
	<ul style="list-style-type: none"> <li>• A kitchen with cookies in the oven baking.</li> <li>• A white oven with cookies being baked inside.</li> <li>• an oven with a pan of cookies baking inside it.</li> <li>• The cookies are inside an oven in the kitchen.</li> </ul>
	<ul style="list-style-type: none"> <li>• Foreign stop sign, possibly in Sanskrit or Cambodian script, with nice tree and water background, off-white property wall typical of India or southeast Asia.</li> <li>• A stop sign in a foreign language by a body of water.</li> <li>• A stop sign posted in a foreign language.</li> <li>• A stop sign has a language that is not English.</li> </ul>

Table 10. Examples of the top-4 captions retrieved for a given image by the proposed jWAE-MSE. Captions that do not match the ground truth are shown in bold.

## Phrase localization

The encoder of the image pipeline passes the 4096-dimensional inputs through two fully connected layers, with 2048 hidden nodes and a 512-dimensional latent space. For the text pipeline, the 6000-dimensional input features are also passed through two fully connected layers, again with 2048 hidden nodes and 512 latent dimensions. The decoders are symmetric to the encoders. We apply ReLU nonlinearities after the hidden layer in the encoder and the decoder. For the discriminator network, we set the size of the three fully connected layers to 384, 256, and 2, respectively. Leaky ReLU nonlinearities are applied to the outputs of first and second layer. The last layer is a linear layer. In Fig. 6, we show CITE [32] built upon our jWAE framework. The inputs to the image pipeline are 4101 dimensional with 4096-dimensional VGG-19 features and 5 spatial location features. Gaussian constraints are applied to the features

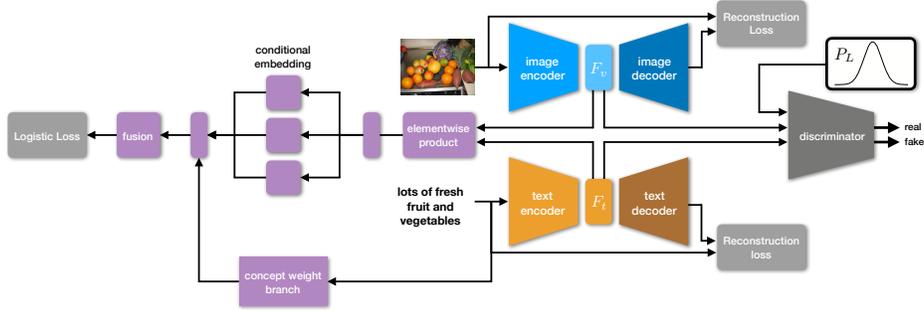


Figure 6. Architecture of CITE [32] with our joint Wasserstein autoencoder.

from image and text pipeline before applying an elementwise product.

### Appendix C. Cross-Modal Retrieval Results

We provide some examples for the cases where our method does not retrieve the correct image for text-to-image retrieval (*c.f.* Table 9). We observe that the retrieved images have similar concepts as the ground truth image. Without requiring any supervised information on image-image similarity, our jWAE is able to align similar concepts in images. This also shows that apart from the ground truth captions of the images, a caption can explain many images in the dataset. In Table 10 we include results for image-to-text retrieval for jWAE-MSE showing that our method is able to recover relevant captions for a given image.

**Generalization across datasets.** In Table 11 and Table 12, we show additional examples for a model trained on COCO and tested on the Flickr30k dataset. We observe that our jWAE retrieves semantically more relevant images (texts) for a given text (image) compared to the competing method. This shows that for good generalization performance of the embeddings, semantic continuity in the latent representations is desirable.

### Appendix D. Phrase Localization Results

In Table 13, we include results of the bounding boxes retrieved by the proposed jWAE for the given phrase and image. Our method is able to retrieve bounding boxes having high overlap with the ground truth bounding box.

Text	Matched images with supervised method [43]	Retrieved image with jWAE-MSE
Two race cars are going down a racetrack bend.		
Two women playing beach volleyball are jumping in the air for the ball.		
Several cyclists make their way through a rutted field, as a camera-man films.		
A hockey goalie stands in front of the net at the ready.		
A woman in white is holding a blue umbrella and walking in the rain.		
Spray obscures the face and torso of the paddle wielding man in the red kayak.		
A line of people are biking down a road.		

Table 11. Examples of images retrieved for a given sentence on Flickr30k based on embeddings trained on COCO.

Image	Supervised approach [43]	Ours (jWAE-MSE)
	<ul style="list-style-type: none"> <li>• People are at a costume party.</li> <li>• A group of parents are sitting at an outdoor assembly.</li> <li>• A group of friends working on a quilt.</li> <li>• A decorative plate with a piece of cake on it.</li> </ul>	<ul style="list-style-type: none"> <li>• A crowd of people in colorful dresses.</li> <li>• Many ethnic people in multicolored robes gather for an event.</li> <li>• A string instrument band and chorus singers rehearsing.</li> <li>• A crowd of women wearing multiple different colors is entering a stone structure.</li> </ul>
	<ul style="list-style-type: none"> <li>• A snowboarder soars through the air.</li> <li>• A stuffed bear with goggles on wearing snow skies.</li> <li>• This is a skier getting some nice jumps in.</li> <li>• A young girl with her bike.</li> </ul>	<ul style="list-style-type: none"> <li>• A snowboarder soars through the air.</li> <li>• A snowboarder is doing a big jump and is in the air.</li> <li>• Somebody is skiing down a mountain.</li> <li>• skier in red pants skiing down a slope.</li> </ul>
	<ul style="list-style-type: none"> <li>• People at a supermarket checkout.</li> <li>• A person wearing a hat made out of yellow bananas.</li> <li>• Two men in orange construction hats are guiding a cart full of brick stones.</li> <li>• A man with two kids looking at pictures on a camera.</li> </ul>	<ul style="list-style-type: none"> <li>• A man is standing in the aisle of a grocery store and staring at the cereal selection.</li> <li>• A customer at the checkout of a grocery store.</li> <li>• A busy store full of people shopping at the store.</li> <li>• A man with a shopping cart is studying the shelves in a supermarket aisle.</li> </ul>
	<ul style="list-style-type: none"> <li>• A soccer player is kicking the ball.</li> <li>• White and red striped bus riding through a city at night.</li> <li>• A man in a blue apron in a kitchen.</li> <li>• A young boy carrying a large soccer ball with a soccer feild in the background.</li> </ul>	<ul style="list-style-type: none"> <li>• A soccer player is running while kicking a ball.</li> <li>• A man in blue plays soccer.</li> <li>• Two girls on separate teams fighting for a soccer ball.</li> <li>• Two girls fight for the soccer ball while playing soccer on the grass.</li> </ul>
	<ul style="list-style-type: none"> <li>• A man jumping down a hill in a forested park.</li> <li>• A man jumps off a hill.</li> <li>• A man sits on a field near a backpack.</li> <li>• A zoo keeper tending to an elephant's mouth.</li> </ul>	<ul style="list-style-type: none"> <li>• Mountain bike riders on a dirt trail.</li> <li>• Woman on bicycle riding down dirt trail.</li> <li>• Bicyclist riding on a dirt race course.</li> <li>• Two women riding their bicycles along a dirt road.</li> </ul>

Table 12. Examples of the top-4 captions retrieved for a given image on Flickr30k based on embeddings trained on COCO.

---

**Images with bounding boxes****Phrases (in red) in sentences**

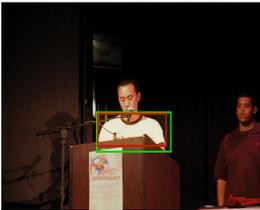
---



A blue, red, and yellow airplane is flying through the air.



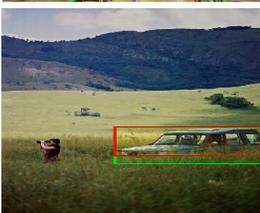
A young man in a t-shirt is speaking at a podium while another young man stands by.



A young man in a t-shirt is speaking at a podium while another young man stands by.



A woman rides her bike by some trees.



A person holding a piece of equipment up to her eyes is standing in a large meadow near a blue vehicle.



A crowd watches as a woman draws caricatures.

---

Table 13. Examples of the bounding boxes retrieved for a phrase in an image by the proposed jWAE. (*green*) Ground truth bounding box; (*red*) Retrieved bounding box.