

---

# Style Transfer from Non-Parallel Text by Cross-Alignment

---

Tianxiao Shen<sup>1</sup> Tao Lei<sup>2</sup> Regina Barzilay<sup>1</sup> Tommi Jaakkola<sup>1</sup>  
<sup>1</sup>MIT CSAIL <sup>2</sup>ASAPP Inc.  
<sup>1</sup>{tianxiao, regina, tommi}@csail.mit.edu <sup>2</sup>tao@asapp.com

## Abstract

This paper focuses on style transfer on the basis of non-parallel text. This is an instance of a broad family of problems including machine translation, decipherment, and sentiment modification. The key challenge is to separate the content from other aspects such as style. We assume a shared latent content distribution across different text corpora, and propose a method that leverages refined alignment of latent representations to perform style transfer. The transferred sentences from one style should match example sentences from the other style as a population. We demonstrate the effectiveness of this cross-alignment method on three tasks: sentiment modification, decipherment of word substitution ciphers, and recovery of word order.<sup>1</sup>

## 1 Introduction

Using massive amounts of parallel data has been essential for recent advances in text generation tasks, such as machine translation and summarization. However, in many text generation problems, we can only assume access to non-parallel or mono-lingual data. Problems such as decipherment or style transfer are all instances of this family of tasks. In all of these problems, we must preserve the content of the source sentence but render the sentence consistent with desired presentation constraints (e.g., style, plaintext/ciphertext).

The goal of controlling one aspect of a sentence such as style independently of its content requires that we can disentangle the two. However, these aspects interact in subtle ways in natural language sentences, and we can succeed in this task only approximately even in the case of parallel data. Our task is more challenging here. We merely assume access to two corpora of sentences with the same distribution of content albeit rendered in different styles. Our goal is to demonstrate that this distributional equivalence of content, if exploited carefully, suffices for us to learn to map a sentence in one style to a style-independent content vector and then decode it to a sentence with the same content but a different style.

In this paper, we introduce a refined alignment of sentence representations across text corpora. We learn an encoder that takes a sentence and its original style indicator as input, and maps it to a style-independent content representation. This is then passed to a style-dependent decoder for rendering. We do not use typical VAEs for this mapping since it is imperative to keep the latent content representation rich and unperturbed. Indeed, richer latent content representations are much harder to align across the corpora and therefore they offer more informative content constraints. Moreover, we reap additional information from cross-generated (style-transferred) sentences, thereby getting two distributional alignment constraints. For example, positive sentences that are style-transferred into negative sentences should match, as a population, the given set of negative sentences. We illustrate this cross-alignment in Figure 1.

---

<sup>1</sup>Our code and data are available at <https://github.com/shentianxiao/language-style-transfer>.

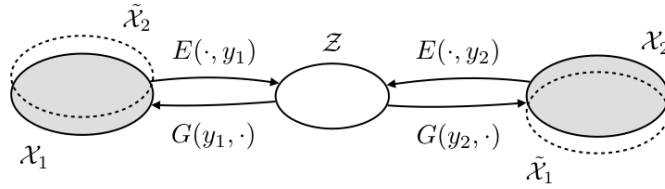


Figure 1: An overview of the proposed cross-alignment method.  $\mathcal{X}_1$  and  $\mathcal{X}_2$  are two sentence domains with different styles  $y_1$  and  $y_2$ , and  $\mathcal{Z}$  is the shared latent content space. Encoder  $E$  maps a sentence to its content representation, and generator  $G$  generates the sentence back when combining with the original style. When combining with a different style, transferred  $\tilde{\mathcal{X}}_1$  is aligned with  $\mathcal{X}_2$  and  $\tilde{\mathcal{X}}_2$  is aligned with  $\mathcal{X}_1$  at the distributional level.

To demonstrate the flexibility of the proposed model, we evaluate it on three tasks: sentiment modification, decipherment of word substitution ciphers, and recovery of word order. In all of these applications, the model is trained on non-parallel data. On the sentiment modification task, the model successfully transfers the sentiment while keeps the content for 41.5% of review sentences according to human evaluation, compared to 41.0% achieved by the control-gen model of Hu et al. (2017). It achieves strong performance on the decipherment and word order recovery tasks, reaching Bleu score of 57.4 and 26.1 respectively, obtaining 50.2 and 20.9 gap than a comparable method without cross-alignment.

## 2 Related work

**Style transfer in vision** Non-parallel style transfer has been extensively studied in computer vision (Gatys et al., 2016; Zhu et al., 2017; Liu and Tuzel, 2016; Liu et al., 2017; Taigman et al., 2016; Kim et al., 2017; Yi et al., 2017). Gatys et al. (2016) explicitly extract content and style features, and then synthesize a new image by combining “content” features of one image with “style” features from another. More recent approaches learn generative networks directly via generative adversarial training (Goodfellow et al., 2014) from two given data domains  $\mathcal{X}_1$  and  $\mathcal{X}_2$ . The key computational challenge in this non-parallel setting is aligning the two domains. For example, CoupledGANs (Liu and Tuzel, 2016) employ weight-sharing between networks to learn cross-domain representation, whereas CycleGAN (Zhu et al., 2017) introduces cycle consistency which relies on transitivity to regularize the transfer functions. While our approach has a similar high-level architecture, the discreteness of natural language does not allow us to reuse these models and necessitates the development of new methods.

**Non-parallel transfer in natural language** In natural language processing, most tasks that involve generation (e.g., translation and summarization) are trained using parallel sentences. Our work most closely relates to approaches that do not utilize parallel data, but instead guide sentence generation from an indirect training signal (Mueller et al., 2017; Hu et al., 2017). For instance, Mueller et al. (2017) manipulate the hidden representation to generate sentences that satisfy a desired property (e.g., sentiment) as measured by a corresponding classifier. However, their model does not necessarily enforce content preservation. More similar to our work, Hu et al. (2017) aims at generating sentences with controllable attributes by learning disentangled latent representations (Chen et al., 2016). Their model builds on variational auto-encoders (VAEs) and uses independency constraints to enforce that attributes can be reliably inferred back from generated sentences. While our model builds on distributional cross-alignment for the purpose of style transfer and content preservation, these constraints can be added in the same way.

**Adversarial training over discrete samples** Recently, a wide range of techniques addresses challenges associated with adversarial training over discrete samples generated by recurrent networks (Yu et al., 2016; Lamb et al., 2016; Hjelm et al., 2017; Che et al., 2017). In our work, we employ the Professor-Forcing algorithm (Lamb et al., 2016) which was originally proposed to close the gap between teacher-forcing during training and self-feeding during testing for recurrent networks. This design fits well with our scenario of style transfer that calls for cross-alignment. By using

continuous relaxation to approximate the discrete sampling process (Jang et al., 2016; Maddison et al., 2016), the training procedure can be effectively optimized through back-propagation (Kusner and Hernández-Lobato, 2016; Goyal et al., 2017).

### 3 Formulation

In this section, we formalize the task of non-parallel style transfer and discuss the feasibility of the learning problem. We assume the data are generated by the following process:

1. a latent style variable  $\mathbf{y}$  is generated from some distribution  $p(\mathbf{y})$ ;
2. a latent content variable  $\mathbf{z}$  is generated from some distribution  $p(\mathbf{z})$ ;
3. a datapoint  $\mathbf{x}$  is generated from conditional distribution  $p(\mathbf{x}|\mathbf{y}, \mathbf{z})$ .

We observe two datasets with the same content distribution but different styles  $\mathbf{y}_1$  and  $\mathbf{y}_2$ , where  $\mathbf{y}_1$  and  $\mathbf{y}_2$  are unknown. Specifically, the two observed datasets  $\mathbf{X}_1 = \{\mathbf{x}_1^{(1)}, \dots, \mathbf{x}_1^{(n)}\}$  and  $\mathbf{X}_2 = \{\mathbf{x}_2^{(1)}, \dots, \mathbf{x}_2^{(m)}\}$  consist of samples drawn from  $p(\mathbf{x}_1|\mathbf{y}_1)$  and  $p(\mathbf{x}_2|\mathbf{y}_2)$  respectively. We want to estimate the style transfer functions between them, namely  $p(\mathbf{x}_1|\mathbf{x}_2; \mathbf{y}_1, \mathbf{y}_2)$  and  $p(\mathbf{x}_2|\mathbf{x}_1; \mathbf{y}_1, \mathbf{y}_2)$ .

A question we must address is when this estimation problem is feasible. Essentially, we only observe the marginal distributions of  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , yet we are going to recover their joint distribution:

$$p(\mathbf{x}_1, \mathbf{x}_2|\mathbf{y}_1, \mathbf{y}_2) = \int_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}_1|\mathbf{y}_1, \mathbf{z})p(\mathbf{x}_2|\mathbf{y}_2, \mathbf{z})d\mathbf{z} \quad (1)$$

As we only observe  $p(\mathbf{x}_1|\mathbf{y}_1)$  and  $p(\mathbf{x}_2|\mathbf{y}_2)$ ,  $\mathbf{y}_1$  and  $\mathbf{y}_2$  are unknown to us. If two different  $\mathbf{y}$  and  $\mathbf{y}'$  lead to the same distribution  $p(\mathbf{x}|\mathbf{y}) = p(\mathbf{x}|\mathbf{y}')$ , then given a dataset  $\mathbf{X}$  sampled from it, its underlying style can be either  $\mathbf{y}$  or  $\mathbf{y}'$ . Consider the following two cases: (1) both datasets  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are sampled from the same style  $\mathbf{y}$ ; (2)  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are sampled from style  $\mathbf{y}$  and  $\mathbf{y}'$  respectively. These two scenarios have different joint distributions, but the observed marginal distributions are the same. To prevent such confusion, we constrain the underlying distributions as stated in the following proposition:

**Proposition 1.** *In the generative framework above,  $\mathbf{x}_1$  and  $\mathbf{x}_2$ 's joint distribution can be recovered from their marginals only if for any different  $\mathbf{y}, \mathbf{y}' \in \mathcal{Y}$ , distributions  $p(\mathbf{x}|\mathbf{y})$  and  $p(\mathbf{x}|\mathbf{y}')$  are different.*

This proposition basically says that  $\mathbf{X}$  generated from different styles should be “distinct” enough, otherwise the transfer task between styles is not well defined. While this seems trivial, it may not hold even for simplified data distributions. The following examples illustrate how the transfer (and recovery) becomes feasible or infeasible under different model assumptions. As we shall see, for a certain family of styles  $\mathcal{Y}$ , the more complex distribution for  $\mathbf{z}$ , the more probable it is to recover the transfer function and the easier it is to search for the transfer.

#### 3.1 Example 1: Gaussian

Consider the common choice that  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  has a centered isotropic Gaussian distribution. Suppose a style  $\mathbf{y} = (\mathbf{A}, \mathbf{b})$  is an affine transformation, i.e.  $\mathbf{x} = \mathbf{A}\mathbf{z} + \mathbf{b} + \epsilon$ , where  $\epsilon$  is a noise variable. For  $\mathbf{b} = \mathbf{0}$  and any orthogonal matrix  $\mathbf{A}$ ,  $\mathbf{A}\mathbf{z} + \mathbf{b} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and hence  $\mathbf{x}$  has the same distribution for any such styles  $\mathbf{y} = (\mathbf{A}, \mathbf{0})$ . In this case, the effect of rotation cannot be recovered.

Interestingly, if  $\mathbf{z}$  has a **more complex distribution**, such as a Gaussian mixture, then affine transformations can be uniquely determined.

**Lemma 1.** *Let  $\mathbf{z}$  be a mixture of Gaussians  $p(\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ . Assume  $K \geq 2$ , and there are two different  $\boldsymbol{\Sigma}_i \neq \boldsymbol{\Sigma}_j$ . Let  $\mathcal{Y} = \{(\mathbf{A}, \mathbf{b}) \mid |\mathbf{A}| \neq 0\}$  be all invertible affine transformations, and  $p(\mathbf{x}|\mathbf{y}, \mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{A}\mathbf{z} + \mathbf{b}, \epsilon^2 \mathbf{I})$ , in which  $\epsilon$  is a noise. Then for all  $\mathbf{y} \neq \mathbf{y}' \in \mathcal{Y}$ ,  $p(\mathbf{x}|\mathbf{y})$  and  $p(\mathbf{x}|\mathbf{y}')$  are different distributions.*

**Theorem 1.** *If the distribution of  $\mathbf{z}$  is a mixture of Gaussians which has more than two different components, and  $\mathbf{x}_1, \mathbf{x}_2$  are two affine transformations of  $\mathbf{z}$ , then the transfer between them can be recovered given their respective marginals.*

### 3.2 Example 2: Word substitution

Consider here another example when  $z$  is a bi-gram language model and a style  $\mathbf{y}$  is a vocabulary in use that maps each “content word” onto its surface form (lexical form). If we observe two realizations  $\mathbf{x}_1$  and  $\mathbf{x}_2$  of the same language  $z$ , the transfer and recovery problem becomes inferring a word alignment between  $\mathbf{x}_1$  and  $\mathbf{x}_2$ .

Note that this is a simplified version of language decipherment or translation. Nevertheless, the recovery problem is still sufficiently hard. To see this, let  $M_1, M_2 \in \mathcal{R}^{n \times n}$  be the estimated bi-gram probability matrix of data  $\mathbf{X}_1$  and  $\mathbf{X}_2$  respectively. Seeking the word alignment is equivalent to finding a permutation matrix  $P$  such that  $P^\top M_1 P \approx M_2$ , which can be expressed as an optimization problem,

$$\min_P \|P^\top M_1 P - M_2\|^2$$

The same formulation applies to graph isomorphism (GI) problems given  $M_1$  and  $M_2$  as the adjacency matrices of two graphs, suggesting that determining the existence and uniqueness of  $P$  is at least GI hard. Fortunately, if  $M$  as a graph is complex enough, the search problem could be more tractable. For instance, if each vertex’s weights of incident edges as a set is unique, then finding the isomorphism can be done by simply matching the sets of edges. This assumption largely applies to our scenario where  $z$  is a complex language model. We empirically demonstrate this in the results section.

The above examples suggest that  $z$  as the latent content variable should carry most complexity of data  $\mathbf{x}$ , while  $\mathbf{y}$  as the latent style variable should have relatively simple effects. We construct the model accordingly in the next section.

## 4 Method

Learning the style transfer function under our generative assumption is essentially learning the conditional distribution  $p(\mathbf{x}_1|\mathbf{x}_2; \mathbf{y}_1, \mathbf{y}_2)$  and  $p(\mathbf{x}_2|\mathbf{x}_1; \mathbf{y}_1, \mathbf{y}_2)$ . Unlike in vision where images are continuous and hence the transfer functions can be learned and optimized directly, the discreteness of language requires us to operate through the latent space. Since  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are conditionally independent given the latent content variable  $z$ ,

$$\begin{aligned} p(\mathbf{x}_1|\mathbf{x}_2; \mathbf{y}_1, \mathbf{y}_2) &= \int_z p(\mathbf{x}_1, z|\mathbf{x}_2; \mathbf{y}_1, \mathbf{y}_2) dz \\ &= \int_z p(z|\mathbf{x}_2, \mathbf{y}_2) \cdot p(\mathbf{x}_1|\mathbf{y}_1, z) dz \\ &= \mathbb{E}_{z \sim p(z|\mathbf{x}_2, \mathbf{y}_2)} [p(\mathbf{x}_1|\mathbf{y}_1, z)] \end{aligned} \quad (2)$$

This suggests us learning an auto-encoder model. Specifically, a style transfer from  $\mathbf{x}_2$  to  $\mathbf{x}_1$  involves two steps—an encoding step that infers  $\mathbf{x}_2$ ’s content  $z \sim p(z|\mathbf{x}_2, \mathbf{y}_2)$ , and a decoding step which generates the transferred counterpart from  $p(\mathbf{x}_1|\mathbf{y}_1, z)$ . In this work, we approximate and train  $p(z|\mathbf{x}, \mathbf{y})$  and  $p(\mathbf{x}|\mathbf{y}, z)$  using neural networks (where  $\mathbf{y} \in \{\mathbf{y}_1, \mathbf{y}_2\}$ ).

Let  $E : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{Z}$  be an encoder that infers the content  $z$  for a given sentence  $\mathbf{x}$  and a style  $\mathbf{y}$ , and  $G : \mathcal{Y} \times \mathcal{Z} \rightarrow \mathcal{X}$  be a generator that generates a sentence  $\mathbf{x}$  from a given style  $\mathbf{y}$  and content  $z$ .  $E$  and  $G$  form an auto-encoder when applying to the same style, and thus we have reconstruction loss,

$$\mathcal{L}_{\text{rec}}(\theta_E, \theta_G) = \mathbb{E}_{\mathbf{x}_1 \sim \mathbf{X}_1} [-\log p_G(\mathbf{x}_1|\mathbf{y}_1, E(\mathbf{x}_1, \mathbf{y}_1))] + \mathbb{E}_{\mathbf{x}_2 \sim \mathbf{X}_2} [-\log p_G(\mathbf{x}_2|\mathbf{y}_2, E(\mathbf{x}_2, \mathbf{y}_2))] \quad (3)$$

where  $\theta$  are the parameters to estimate.

In order to make a meaningful transfer by flipping the style,  $\mathbf{X}_1$  and  $\mathbf{X}_2$ ’s content space must coincide, as our generative framework presumed. To constrain that  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are generated from the same latent content distribution  $p(z)$ , one option is to apply a variational auto-encoder (Kingma and Welling, 2013). A VAE imposes a prior density  $p(z)$ , such as  $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and uses a KL-divergence regularizer to align both posteriors  $p_E(z|\mathbf{x}_1, \mathbf{y}_1)$  and  $p_E(z|\mathbf{x}_2, \mathbf{y}_2)$  to it,

$$\mathcal{L}_{\text{KL}}(\theta_E) = \mathbb{E}_{\mathbf{x}_1 \sim \mathbf{X}_1} [D_{\text{KL}}(p_E(z|\mathbf{x}_1, \mathbf{y}_1)||p(z))] + \mathbb{E}_{\mathbf{x}_2 \sim \mathbf{X}_2} [D_{\text{KL}}(p_E(z|\mathbf{x}_2, \mathbf{y}_2)||p(z))] \quad (4)$$

The overall objective is to minimize  $\mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{KL}}$ , whose opposite is the variational lower bound of data likelihood.

However, as we have argued in the previous section, restricting  $z$  to a simple and even distribution and pushing most complexity to the decoder may not be a good strategy for non-parallel style transfer. In contrast, a standard auto-encoder simply minimizes the reconstruction error, encouraging  $z$  to carry as much information about  $x$  as possible. On the other hand, it lowers the entropy in  $p(x|y, z)$ , which helps to produce meaningful style transfer in practice as we flip between  $y_1$  and  $y_2$ . Without explicitly modeling  $p(z)$ , it is still possible to force distributional alignment of  $p(z|y_1)$  and  $p(z|y_2)$ . To this end, we introduce two constrained variants of auto-encoder.

#### 4.1 Aligned auto-encoder

Dispense with VAEs that make an explicit assumption about  $p(z)$  and align both posteriors to it, we align  $p_E(z|y_1)$  and  $p_E(z|y_2)$  with each other, which leads to the following constrained optimization problem:

$$\begin{aligned} \theta^* &= \arg \min_{\theta} \mathcal{L}_{\text{rec}}(\theta_E, \theta_G) \\ \text{s.t. } E(x_1, y_1) &\stackrel{d}{=} E(x_2, y_2) \quad x_1 \sim X_1, x_2 \sim X_2 \end{aligned} \quad (5)$$

In practice, a Lagrangian relaxation of the primal problem is instead optimized. We introduce an adversarial discriminator  $D$  to align the aggregated posterior distribution of  $z$  from different styles (Makhzani et al., 2015).  $D$  aims to distinguish between these two distributions:

$$\mathcal{L}_{\text{adv}}(\theta_E, \theta_D) = \mathbb{E}_{x_1 \sim X_1}[-\log D(E(x_1, y_1))] + \mathbb{E}_{x_2 \sim X_2}[-\log(1 - D(E(x_2, y_2)))] \quad (6)$$

The overall training objective is a min-max game played among the encoder  $E$ , generator  $G$  and discriminator  $D$ . They constitute an aligned auto-encoder:

$$\min_{E, G} \max_D \mathcal{L}_{\text{rec}} - \lambda \mathcal{L}_{\text{adv}} \quad (7)$$

We implement the encoder  $E$  and generator  $G$  using single-layer RNNs with GRU cell.  $E$  takes an input sentence  $x$  with initial hidden state  $y$ , and outputs the last hidden state  $z$  as its content representation.  $G$  generates a sentence  $x$  conditioned on latent state  $(y, z)$ . To align the distributions of  $z_1 = E(x_1, y_1)$  and  $z_2 = E(x_2, y_2)$ , the discriminator  $D$  is a feed-forward network with a single hidden layer and a sigmoid output layer.

#### 4.2 Cross-aligned auto-encoder

The second variant, cross-aligned auto-encoder, directly aligns the transferred samples from one style with the true samples from the other. Under the generative assumption,  $p(x_2|y_2) = \int_{x_1} p(x_2|x_1; y_1, y_2)p(x_1|y_1)dx_1$ , thus  $x_2$  (sampled from the left-hand side) should exhibit the same distribution as transferred  $x_1$  (sampled from the right-hand side), and vice versa. Similar to our first model, the second model uses two discriminators  $D_1$  and  $D_2$  to align the populations.  $D_1$ 's job is to distinguish between real  $x_1$  and transferred  $x_2$ , and  $D_2$ 's job is to distinguish between real  $x_2$  and transferred  $x_1$ .

Adversarial training over the discrete samples generated by  $G$  hinders gradients propagation. Although sampling-based gradient estimator such as REINFORCE (Williams, 1992) can be adopted, training with these methods can be unstable due to the high variance of the sampled gradient. Instead, we employ two recent techniques to approximate the discrete training (Hu et al., 2017; Lamb et al., 2016). First, instead of feeding a single sampled word as the input to the generator RNN, we use the softmax distribution over words instead. Specifically, during the generating process of transferred  $x_2$  from  $G(y_1, z_2)$ , suppose at time step  $t$  the output logit vector is  $v_t$ . We feed its peaked distribution  $\text{softmax}(v_t/\gamma)$  as the next input, where  $\gamma \in (0, 1)$  is a temperature parameter.

Secondly, we use Professor-Forcing (Lamb et al., 2016) to match the sequence of hidden states instead of the output words, which contains the information about outputs and is smoothly distributed. That is, the input to the discriminator  $D_1$  is the sequence of hidden states of either (1)  $G(y_1, z_1)$  teacher-forced by a real example  $x_1$ , or (2)  $G(y_1, z_2)$  self-fed by previous soft distributions.

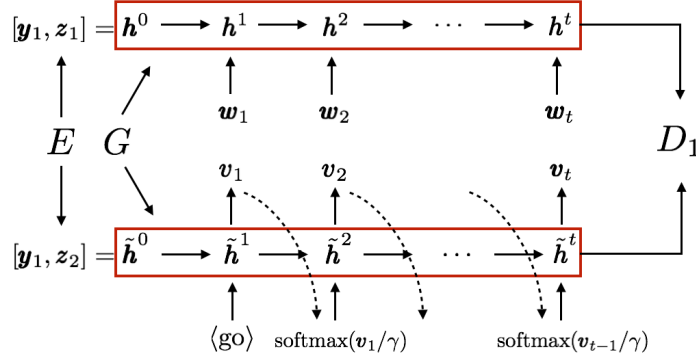


Figure 2: Cross-aligning between  $x_1$  and transferred  $x_2$ . For  $x_1$ ,  $G$  is teacher-forced by its words  $w_1 w_2 \dots w_t$ . For transferred  $x_2$ ,  $G$  is self-fed by previous output logits. The sequence of hidden states  $h^0, \dots, h^t$  and  $\tilde{h}^0, \dots, \tilde{h}^t$  are passed to discriminator  $D_1$  to be aligned. Note that our first variant aligned auto-encoder is a special case of this, where only  $h^0$  and  $\tilde{h}^0$ , i.e.  $z_1$  and  $z_2$ , are aligned.

---

**Algorithm 1** Cross-aligned auto-encoder training. The hyper-parameters are set as  $\lambda = 1, \gamma = 0.001$  and learning rate is 0.0001 for all experiments in this paper.

---

**Input:** Two corpora of different styles  $X_1, X_2$ . Lagrange multiplier  $\lambda$ , temperature  $\gamma$ .

Initialize  $\theta_E, \theta_G, \theta_{D_1}, \theta_{D_2}$

**repeat**

**for**  $p = 1, 2; q = 2, 1$  **do**

Sample a mini-batch of  $k$  examples  $\{x_p^{(i)}\}_{i=1}^k$  from  $X_p$

Get the latent content representations  $z_p^{(i)} = E(x_p^{(i)}, y_p)$

Unroll  $G$  from initial state  $(y_p, z_p^{(i)})$  by feeding  $x_p^{(i)}$ , and get the hidden states sequence  $h_p^{(i)}$

Unroll  $G$  from initial state  $(y_q, z_p^{(i)})$  by feeding previous soft output distribution with temperature  $\gamma$ , and get the transferred hidden states sequence  $\tilde{h}_p^{(i)}$

**end for**

Compute the reconstruction  $\mathcal{L}_{\text{rec}}$  by Eq. (3)

Compute  $D_1$ 's (and symmetrically  $D_2$ 's) loss:

$$\mathcal{L}_{\text{adv}_1} = -\frac{1}{k} \sum_{i=1}^k \log D_1(h_1^{(i)}) - \frac{1}{k} \sum_{i=1}^k \log(1 - D_1(\tilde{h}_2^{(i)})) \quad (8)$$

Update  $\{\theta_E, \theta_G\}$  by gradient descent on loss

$$\mathcal{L}_{\text{rec}} - \lambda(\mathcal{L}_{\text{adv}_1} + \mathcal{L}_{\text{adv}_2}) \quad (9)$$

Update  $\theta_{D_1}$  and  $\theta_{D_2}$  by gradient descent on loss  $\mathcal{L}_{\text{adv}_1}$  and  $\mathcal{L}_{\text{adv}_2}$  respectively

**until** convergence

**Output:** Style transfer functions  $G(y_2, E(\cdot, y_1)) : \mathcal{X}_1 \rightarrow \mathcal{X}_2$  and  $G(y_1, E(\cdot, y_2)) : \mathcal{X}_2 \rightarrow \mathcal{X}_1$

---

The running procedure of our cross-aligned auto-encoder is illustrated in Figure 2. Note that cross-aligning strengthens the alignment of latent variable  $z$  over the recurrent network of generator  $G$ . By aligning the whole sequence of hidden states, it prevents  $z_1$  and  $z_2$ 's initial misalignment from propagating through the recurrent generating process, as a result of which the transferred sentence may end up somewhere far from the target domain.

We implement both  $D_1$  and  $D_2$  using convolutional neural networks for sequence classification (Kim, 2014). The training algorithm is presented in Algorithm 1.

## 5 Experimental setup

**Sentiment modification** Our first experiment focuses on text rewriting with the goal of changing the underlying sentiment, which can be regarded as “style transfer” between negative and positive sentences. We run experiments on Yelp restaurant reviews, utilizing readily available user ratings associated with each review. Following standard practice, reviews with rating above three are considered positive, and those below three are considered negative. While our model operates at the sentence level, the sentiment annotations in our dataset are provided at the document level. We assume that all the sentences in a document have the same sentiment. This is clearly an oversimplification, since some sentences (e.g., background) are sentiment neutral. Given that such sentences are more common in long reviews, we filter out reviews that exceed 10 sentences. We further filter the remaining sentences by eliminating those that exceed 15 words. The resulting dataset has 250K negative sentences, and 350K positive ones. The vocabulary size is 10K after replacing words occurring less than 5 times with the “<unk>” token. As a baseline model, we compare against the control-gen model of Hu et al. (2017).

To quantitatively evaluate the transferred sentences, we adopt a model-based evaluation metric similar to the one used for image transfer (Isola et al., 2016). Specifically, we measure how often a transferred sentence has the correct sentiment according to a pre-trained sentiment classifier. For this purpose, we use the TextCNN model as described in Kim (2014). On our simplified dataset for style transfer, it achieves nearly perfect accuracy of 97.4%.

While the quantitative evaluation provides some indication of transfer quality, it does not capture all the aspects of this generation task. Therefore, we also perform two human evaluations on 500 sentences randomly selected from the test set<sup>2</sup>. In the first evaluation, the judges were asked to rank generated sentences in terms of their fluency and sentiment. Fluency was rated from 1 (unreadable) to 4 (perfect), while sentiment categories were “positive”, “negative”, or “neither” (which could be contradictory, neutral or nonsensical). In the second evaluation, we evaluate the transfer process comparatively. The annotator was shown a source sentence and the corresponding outputs of the systems in a random order, and was asked “Which transferred sentence is semantically equivalent to the source sentence with an opposite sentiment?”. They can be both satisfactory, A/B is better, or both unsatisfactory. We collect two labels for each question. The label agreement and conflict resolution strategy can be found in the supplementary material. Note that the two evaluations are not redundant. For instance, a system that always generates the same grammatically correct sentence with the right sentiment independently of the source sentence will score high in the first evaluation setup, but low in the second one.

**Word substitution decipherment** Our second set of experiments involves decipherment of word substitution ciphers, which has been previously explored in NLP literature (Dou and Knight, 2012; Nuhn and Ney, 2013). These ciphers replace every word in plaintext (natural language) with a cipher token according to a 1-to-1 substitution key. The decipherment task is to recover the plaintext from ciphertext. It is trivial if we have access to parallel data. However we are interested to consider a non-parallel decipherment scenario. For training, we select 200K sentences as  $X_1$ , and apply a substitution cipher  $f$  on a different set of 200K sentences to get  $X_2$ . While these sentences are non-parallel, they are drawn from the same distribution from the review dataset. The development and test sets have 100K parallel sentences  $D_1 = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$  and  $D_2 = \{f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(n)})\}$ . We can quantitatively compare between  $D_1$  and transferred (deciphered)  $D_2$  using Bleu score (Papineni et al., 2002).

Clearly, the difficulty of this decipherment task depends on the number of substituted words. Therefore, we report model performance with respect to the percentage of the substituted vocabulary. Note that the transfer models do not know that  $f$  is a word substitution function. They learn it entirely from the data distribution.

In addition to having different transfer models, we introduce a simple decipherment baseline based on word frequency. Specifically, we assume that words shared between  $X_1$  and  $X_2$  do not require translation. The rest of the words are mapped based on their frequency, and ties are broken arbitrarily. Finally, to assess the difficulty of the task, we report the accuracy of a machine translation system trained on a parallel corpus (Klein et al., 2017).

<sup>2</sup>we eliminated 37 sentences from them that were judged as neutral by human judges.

Method	accuracy
Hu et al. (2017)	83.5
Variational auto-encoder	23.2
Aligned auto-encoder	48.3
Cross-aligned auto-encoder	78.4

Table 1: Sentiment accuracy of transferred sentences, as measured by a pretrained classifier.

Method	sentiment	fluency	overall transfer
Hu et al. (2017)	70.8	3.2	41.0
Cross-align	62.6	2.8	41.5

Table 2: Human evaluations on sentiment, fluency and overall transfer quality. Fluency rating is from 1 (unreadable) to 4 (perfect). Overall transfer quality is evaluated in a comparative manner, where the judge is shown a source sentence and two transferred sentences, and decides whether they are both good, both bad, or one is better.

**Word order recovery** Our final experiments focus on the word ordering task, also known as bag translation (Brown et al., 1990; Schmaltz et al., 2016). By learning the style transfer functions between original English sentences  $X_1$  and shuffled English sentences  $X_2$ , the model can be used to recover the original word order of a shuffled sentence (or conversely to randomly permute a sentence). The process to construct non-parallel training data and parallel testing data is the same as in the word substitution decipherment experiment. Again the transfer models do not know that  $f$  is a shuffle function and learn it completely from data.

## 6 Results

**Sentiment modification** Table 1 and Table 2 show the performance of various models for both human and automatic evaluation. The control-gen model of Hu et al. (2017) performs better in terms of sentiment accuracy in both evaluations. This is not surprising as their generation is directly guided by a sentiment classifier. Their system also achieves higher fluency score. However, these gains do not translate into improvements in terms of the overall transfer, where our model fared better. As can be seen from the examples listed in Table 3, our model is more consistent with the grammatical structure and semantic meaning of the source sentence. In contrast, their model achieves sentiment change by generating an entirely new sentence which has little overlap with the source. The discrepancy between the two experiments demonstrates the crucial importance of developing appropriate evaluation measures to compare models for style transfer.

**Word substitution decipherment** Table 4 summarizes the performance of our model and the baselines on the decipherment task, at various levels of word substitution. Consistent with our intuition, the last row in this table shows that the task is trivial when the parallel data is provided. In non-parallel case, the difficulty of the task is driven by the substitution rate. Across all the testing conditions, our cross-aligned model consistently outperforms its counterparts. The difference becomes more pronounced as the task becomes harder. When the substitution rate is 20%, all methods do a reasonably good job in recovering substitutions. However, when 100% of the words are substituted (as expected in real language decipherment), the poor performance of variational autoencoder and aligned auto-encoder rules out their application for this task.

**Word order recovery** The last column in Table 4 demonstrates the performance on the word order recovery task. Order recovery is much harder—even when trained with parallel data, the machine translation model achieves only 64.6 Bleu score. Note that some generated orderings may be completely valid (e.g., reordering conjunctions), but the models will be penalized for producing them. In this task, only the cross-aligned auto-encoder achieves grammatical reorder to a certain extent, demonstrated by its Bleu score 26.1. Other models fail this task, doing no better than no transfer.



From negative to positive
consistently slow . consistently good . consistently fast .
my goodness it was so gross . my husband 's steak was phenomenal . my goodness was so awesome .
it was super dry and had a weird taste to the entire slice . it was a great meal and the tacos were very kind of good . it was super flavorful and had a nice texture of the whole side .
From positive to negative
i love the ladies here ! i avoid all the time ! i hate the doctor here !
my appetizer was also very good and unique . my bf was n't too pleased with the beans . my appetizer was also very cold and not fresh whatsoever .
came here with my wife and her grandmother ! came here with my wife and hated her ! came here with my wife and her son .

Table 3: Sentiment transfer samples. The first line is an input sentence, the second and third lines are the generated sentences after sentiment transfer by Hu et al. (2017) and our cross-aligned auto-encoder, respectively.

Method	Substitution decipher					Order recover
	20%	40%	60%	80%	100%	
No transfer (copy)	56.4	21.4	6.3	4.5	0	5.1
Unigram matching	74.3	48.1	17.8	10.7	1.2	-
Variational auto-encoder	79.8	59.6	44.6	34.4	0.9	5.3
Aligned auto-encoder	81.0	68.9	50.7	45.6	7.2	5.2
Cross-aligned auto-encoder	<b>83.8</b>	<b>79.1</b>	<b>74.7</b>	<b>66.1</b>	<b>57.4</b>	<b>26.1</b>
Parallel translation	99.0	98.9	98.2	98.5	97.2	64.6

Table 4: Bleu scores of word substitution decipherment and word order recovery.

## 7 Conclusion

Transferring languages from one style to another has been previously trained using parallel data. In this work, we formulate the task as a *decipherment problem* with access only to non-parallel data. The two data collections are assumed to be generated by a latent variable generative model. Through this view, our method optimizes neural networks by forcing distributional alignment (invariance) over the latent space or sentence populations. We demonstrate the effectiveness of our method on tasks that permit quantitative evaluation, such as sentiment transfer, word substitution decipherment and word ordering. The decipherment view also provides an interesting open question—*when can the joint distribution  $p(\mathbf{x}_1, \mathbf{x}_2)$  be recovered given only marginal distributions?* We believe addressing this general question would promote the style transfer research in both vision and NLP.

## Acknowledgments

We thank Nicholas Matthews for helping to facilitate human evaluations, and Zhiting Hu for sharing his code. We also thank Jonas Mueller, Arjun Majumdar, Olga Simek, Danelle Shah, MIT NLP group and the reviewers for their helpful comments. This work was supported by MIT Lincoln Laboratory.

## References

- Peter F Brown, John Cocke, Stephen A Della Pietra, Vincent J Della Pietra, Fredrick Jelinek, John D Lafferty, Robert L Mercer, and Paul S Roossin. A statistical approach to machine translation. *Computational linguistics*, 16(2):79–85, 1990.
- Tong Che, Yanran Li, Ruixiang Zhang, R Devon Hjelm, Wenjie Li, Yangqiu Song, and Yoshua Bengio. Maximum-likelihood augmented discrete generative adversarial networks. *arXiv preprint arXiv:1702.07983*, 2017.
- Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in neural information processing systems*, 2016.
- Qing Dou and Kevin Knight. Large scale decipherment for out-of-domain machine translation. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 266–275. Association for Computational Linguistics, 2012.
- Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2414–2423, 2016.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- Kartik Goyal, Chris Dyer, and Taylor Berg-Kirkpatrick. Differentiable scheduled sampling for credit assignment. *arXiv preprint arXiv:1704.06970*, 2017.
- R Devon Hjelm, Athul Paul Jacob, Tong Che, Kyunghyun Cho, and Yoshua Bengio. Boundary-seeking generative adversarial networks. *arXiv preprint arXiv:1702.08431*, 2017.
- Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P Xing. Controllable text generation. *arXiv preprint arXiv:1703.00955*, 2017.
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. *arXiv preprint arXiv:1611.07004*, 2016.
- Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- Taeksoo Kim, Moonsoo Cha, Hyunsoo Kim, Jungkwon Lee, and Jiwon Kim. Learning to discover cross-domain relations with generative adversarial networks. *arXiv preprint arXiv:1703.05192*, 2017.
- Yoon Kim. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*, 2014.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander M Rush. Opennmt: Open-source toolkit for neural machine translation. *arXiv preprint arXiv:1701.02810*, 2017.
- Matt J Kusner and José Miguel Hernández-Lobato. Gans for sequences of discrete elements with the gumbel-softmax distribution. *arXiv preprint arXiv:1611.04051*, 2016.

- Alex M Lamb, Anirudh Goyal ALIAS PARTH GOYAL, Ying Zhang, Saizheng Zhang, Aaron C Courville, and Yoshua Bengio. Professor forcing: A new algorithm for training recurrent networks. In *Advances In Neural Information Processing Systems*, pages 4601–4609, 2016.
- Ming-Yu Liu and Oncel Tuzel. Coupled generative adversarial networks. In *Advances in Neural Information Processing Systems*, pages 469–477, 2016.
- Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. *arXiv preprint arXiv:1703.00848*, 2017.
- Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015.
- Jonas Mueller, Tommi Jaakkola, and David Gifford. Sequence to better sequence: continuous revision of combinatorial structures. *International Conference on Machine Learning (ICML)*, 2017.
- Malte Nuhn and Hermann Ney. Decipherment complexity in 1: 1 substitution ciphers. In *ACL (1)*, pages 615–621, 2013.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics, 2002.
- Allen Schmaltz, Alexander M. Rush, and Stuart Shieber. Word ordering without syntax. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2319–2324. Association for Computational Linguistics, 2016.
- Yaniv Taigman, Adam Polyak, and Lior Wolf. Unsupervised cross-domain image generation. *arXiv preprint arXiv:1611.02200*, 2016.
- Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- Zili Yi, Hao Zhang, Ping Tan Gong, et al. Dualgan: Unsupervised dual learning for image-to-image translation. *arXiv preprint arXiv:1704.02510*, 2017.
- Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. Seqgan: sequence generative adversarial nets with policy gradient. *arXiv preprint arXiv:1609.05473*, 2016.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv preprint arXiv:1703.10593*, 2017.

## A Proof of Lemma 1

**Lemma 1.** Let  $\mathbf{z}$  be a mixture of Gaussians  $p(\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ . Assume  $K \geq 2$ , and there are two different  $\boldsymbol{\Sigma}_i \neq \boldsymbol{\Sigma}_j$ . Let  $\mathcal{Y} = \{(\mathbf{A}, \mathbf{b}) \mid |\mathbf{A}| \neq 0\}$  be all invertible affine transformations, and  $p(\mathbf{x}|\mathbf{y}, \mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{A}\mathbf{z} + \mathbf{b}, \epsilon^2 \mathbf{I})$ , in which  $\epsilon$  is a noise. Then for all  $\mathbf{y} \neq \mathbf{y}' \in \mathcal{Y}$ ,  $p(\mathbf{x}|\mathbf{y})$  and  $p(\mathbf{x}|\mathbf{y}')$  are different distributions.

*Proof.*

$$p(\mathbf{x}|\mathbf{y} = (\mathbf{A}, \mathbf{b})) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \mathbf{A}\boldsymbol{\mu}_k + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}_k\mathbf{A}^\top + \epsilon^2 \mathbf{I})$$

For different  $\mathbf{y} = (\mathbf{A}, \mathbf{b})$  and  $\mathbf{y}' = (\mathbf{A}', \mathbf{b}')$ ,  $p(\mathbf{x}|\mathbf{y}) = p(\mathbf{x}|\mathbf{y}')$  entails that for  $k = 1, \dots, K$ ,

$$\begin{cases} \mathbf{A}\boldsymbol{\mu}_k + \mathbf{b} = \mathbf{A}'\boldsymbol{\mu}_k + \mathbf{b}' \\ \mathbf{A}\boldsymbol{\Sigma}_k\mathbf{A}^\top = \mathbf{A}'\boldsymbol{\Sigma}_k\mathbf{A}'^\top \end{cases}$$

Since all  $\mathcal{Y}$  are invertible,

$$(\mathbf{A}^{-1}\mathbf{A}')\boldsymbol{\Sigma}_k(\mathbf{A}^{-1}\mathbf{A}')^\top = \boldsymbol{\Sigma}_k$$

Suppose  $\boldsymbol{\Sigma}_k = \mathbf{Q}_k\mathbf{D}_k\mathbf{Q}_k^\top$  is  $\boldsymbol{\Sigma}_k$ 's orthogonal diagonalization. If  $k = 1$ , all solutions for  $\mathbf{A}^{-1}\mathbf{A}'$  have the form:

$$\left\{ \mathbf{Q}\mathbf{D}^{1/2}\mathbf{U}\mathbf{D}^{-1/2}\mathbf{Q}^\top \mid \mathbf{U} \text{ is orthogonal} \right\}$$

However, when  $K \geq 2$  and there are two different  $\boldsymbol{\Sigma}_i \neq \boldsymbol{\Sigma}_j$ , the only solution is  $\mathbf{A}^{-1}\mathbf{A}' = \mathbf{I}$ , i.e.  $\mathbf{A} = \mathbf{A}'$ , and thus  $\mathbf{b} = \mathbf{b}'$ .

Therefore, for all  $\mathbf{y} \neq \mathbf{y}'$ ,  $p(\mathbf{x}|\mathbf{y}) \neq p(\mathbf{x}|\mathbf{y}')$ . □