

# SentenceMIM: A Latent Variable Language Model

Micha Livne<sup>1,2</sup> Kevin Swersky<sup>3</sup> David J. Fleet<sup>1,2,3</sup>

## Abstract

SentenceMIM is a probabilistic auto-encoder for language data, trained with Mutual Information Machine (MIM) learning to provide a fixed length representation of variable length language observations (*i.e.*, similar to VAE). Previous attempts to learn VAEs for language data faced challenges due to posterior collapse. MIM learning encourages high mutual information between observations and latent variables, and is robust against posterior collapse. As such, it learns informative representations whose dimension can be an order of magnitude higher than existing language VAEs. Importantly, the SentenceMIM loss has no hyper-parameters, simplifying optimization. We compare sentenceMIM with VAE, and AE on multiple datasets. SentenceMIM yields excellent reconstruction, comparable to AEs, with a rich structured latent space, comparable to VAEs. The structured latent representation is demonstrated with interpolation between sentences of different lengths. We demonstrate the versatility of sentenceMIM by utilizing a trained model for question-answering and transfer learning, without fine-tuning, outperforming VAE and AE with similar architectures.

## 1. Introduction

Generative modelling of text has become one of the predominant approaches to natural language processing (NLP), particularly in the machine learning community. It is favoured because it supports probabilistic reasoning and it provides a principled framework for unsupervised learning in the form of maximum likelihood. Unlike computer vision, where various generative approaches have proliferated (Dinh et al., 2017; Goodfellow et al., 2014; Kingma & Welling, 2013; Oord et al., 2016; Rezende et al., 2014; Vahdat &

Kautz, 2020), current methods for text mainly rely on auto-regressive models (*e.g.*, Brown et al. (2020)).

Generative latent variable models (LVMs), such as the variational auto-encoder (VAE) (Kingma & Welling, 2013; Rezende et al., 2014), are versatile and have been successfully applied to myriad domains. Such models consist of an encoder, which maps observations to distributions over latent codes, and a decoder that maps latent codes to distributions over observations. LVMs are widely used and studied because they can learn a latent representation that carries many useful properties. Observations are encoded as fixed-length vectors that capture salient information, allowing for semantic comparison, interpolation, and search. They are often useful in support of downstream tasks, such as transfer or k-shot learning. They are also often interpretable, capturing distinct factors of variation in different latent dimensions. These properties have made LVMs especially compelling in the vision community.

Despite their desirable qualities, generative LVMs have not enjoyed the same level of success with text data. There have been recent proposals to adapt VAEs to text (Bowman et al., 2015; Guu et al., 2017; Kruengkrai, 2019; Li et al., 2019b; Yang et al., 2017; Bosc & Vincent, 2020), but despite encouraging progress, they have not reached the same level of performance on natural language benchmarks as auto-regressive models (*e.g.*, (Merity et al., 2017; Rae et al., 2018; Wang et al., 2019)). This is often attributed to the phenomenon of posterior collapse (Fang et al., 2019; Li et al., 2019a), in which the decoder captures all of the modelling power and the encoder conveys little to no information. For text, where the decoder is naturally auto-regressive, this has proven challenging to mitigate. A notable exception by Li et al. (2020) utilizes pre-trained BERT encoder (Devlin et al., 2019) and GPT-2 decoder (Radford et al., 2019) in order to train a powerful VAE model. While showing strong PPL results, the training requires carefully designed heuristics to reduce posterior collapse.

This paper introduces sentenceMIM (sMIM), a new LVM for text. We use the Mutual Information Machine (MIM) framework by Livne et al. (2019) for learning, and base our architecture on Bowman et al. (2015). MIM is a recently introduced LVM framework that shares the same underlying architecture as VAEs, but uses a different learning objective

<sup>1</sup>Department of Computer Science, University of Toronto

<sup>2</sup>Vector Institute <sup>3</sup>Google Research. Correspondence to: Micha Livne <mlivne@cs.toronto.edu>, Kevin Swersky <kswersky@google.com>, David J. Fleet <fleet@cs.toronto.edu>.

that is robust against posterior collapse. MIM learns a highly informative and compressed latent representation, and often strictly benefits from more powerful architectures.

We argue that an ideal LVM should be able to capture all aspects of variation of variable-size text observations within a fixed-size latent representation. As such, high-dimensional latent codes are required, which is challenging with VAEs. An ideal model should provide excellent reconstruction, with fixed-size codes for variable-length sentences, and be useful for various downstream tasks.

Auto-encoders (AEs) (Hinton & Zemel, 1994) provide excellent reconstruction, but lack useful semantic structure in the learned representations. AEs also allow one to learn high dimensional latent codes, only limited in practice by over-fitting. VAEs, on the other hand, encourage semantic representations by regularizing the distribution over latent codes to match a given prior. Such regularization, however, also contributes to posterior collapse, limiting the dimension of latent codes and reconstruction quality. Here we propose MIM learning to enable high dimensional representations, while encouraging low latent entropy (under certain conditions) to promote clustering of semantically similar observations. By encouraging the latent codes to match a known distribution we preserve the ability generate samples. The resulting model offers a learned representation with high mutual information (*i.e.*, to capture aspects of variation in the data), with low marginal entropy (*i.e.*, introducing semantic structure to the learned representation), while aligning the latent distribution with a known prior. sMIM also requires no hyper-parameter tuning for the loss, similar to AEs, which simplifies training.

This paper explores and contrasts properties of VAE learning, MIM learning, and AE learning on four well-known text datasets, all with similar architectures. We show that sMIM provides better reconstruction than VAE models, matching the reconstruction accuracy of AEs, but with semantically meaningful representations, comparable to VAEs. We further demonstrate the quality of the sMIM representation by generating diverse samples around a given sentence and interpolating between sentences. Finally, we show the versatility of the learned representation by applying a pre-trained sMIM model to a question answering task with state-of-art performance as compared to single task, supervised models.

## 2. Problem Formulation

Let  $\mathbf{x} \in \mathcal{X} = \{\mathbf{x}_i\}_{i=1}^X$  be a discrete variable representing a sequence of  $T$  tokens from a finite vocabulary  $\mathcal{V}$ . A sequence might be a sentence or a paragraph, for example. The set  $\mathcal{X}$  comprises all sequences we aim to model. The size of  $\mathcal{X}$ , *i.e.*,  $X$ , is typically unknown and large. Let  $\mathcal{P}(\mathbf{x})$  be the unknown probability of sentence  $\mathbf{x} \in \mathcal{X}$ .

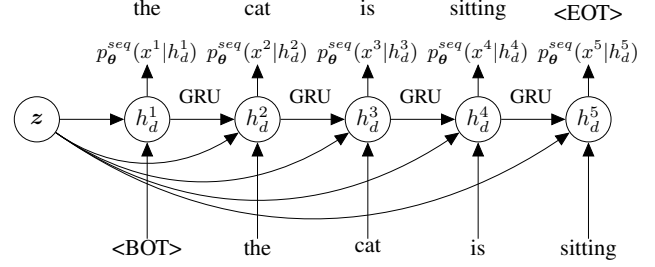


Figure 1. The decoder, implemented with GRU, is auto-regressive and conditioned on latent code  $\mathbf{z}$ . Words are represented by parametric embeddings. At each step the inputs are the latent code and previous output token. The GRU output provides a categorical distribution over tokens  $\mathbf{x}^k$ , from which the next token is sampled.

We model the distribution over a sequence of length  $T$  as an auto-regressive distribution over  $T + 1$  tokens, where the additional end-of-text special token<sup>1</sup>,  $\langle \text{EOT} \rangle$ , effectively captures the probability that the sequence length is  $T$ . More explicitly, we model the following distribution

$$\mathcal{P}(\mathbf{x}) = \sum_{k=0}^T p(\mathbf{x}^k | \mathbf{x}^{<k}) \quad (1)$$

where  $\mathbf{x}^{<k}$  denotes the tokens preceding  $\mathbf{x}^k$ ,  $p(\mathbf{x}^k | \mathbf{x}^{<k})$  is a categorical distribution over  $\mathcal{V}$ , and  $p(\mathbf{x}^T = \langle \text{EOT} \rangle | \mathbf{x}^{<T})$  is the probability that  $T$  is the sentence length.

We learn a latent variable model given  $N$  fair samples from  $\mathcal{P}(\mathbf{x})$ , where  $N \ll X$ , with discrete observations  $\mathbf{x} \in \mathcal{X}$ , and a continuous latent space  $\mathbf{z} \in \mathbb{R}^d$ . The encoder,  $q_\theta(\mathbf{z} | \mathbf{x})$ , maps sequences to a distribution over continuous latent codes, and a corresponding decoder,  $p_\theta(\mathbf{x} | \mathbf{z})$ , maps a latent code to a distribution over sequences. Let  $\theta$  be the joint parameters of the encoder and decoder.

### 2.1. Encoder-Decoder Specification

In what follows we adapt the architecture proposed by Bowman et al. (2015), the main difference being the use of GRUs (Cho et al., 2014) instead of LSTMs (Hochreiter & Schmidhuber, 1997). We opt for a simple architecture instead of more recent variants, such as Transformers (Vaswani et al., 2017) or AWD-LSTMs (Merity et al., 2018), to focus on the effect of the learning framework on a given architecture (*i.e.*, MIM, VAE, AE), rather than the architecture itself.

Beginning with the generative process, let  $p_\theta(\mathbf{x} | \mathbf{z})$  be a conditional auto-regressive distribution over a sequence of  $T$  tokens,  $\mathbf{x} = (\mathbf{x}^0, \dots, \mathbf{x}^{T-1}, \mathbf{x}^T = \langle \text{EOT} \rangle)$ ,

$$\log p_\theta(\mathbf{x} | \mathbf{z}) = \sum_{k=0}^T \log p_\theta(\mathbf{x}^k | \mathbf{x}^{<k}, \mathbf{z}) \quad (2)$$

<sup>1</sup> $\langle \text{BOT} \rangle$ ,  $\langle \text{EOT} \rangle$  are a special beginning/end-of-text tokens. The token  $\langle \text{UNK} \rangle$  represents an out-of-vocabulary word.

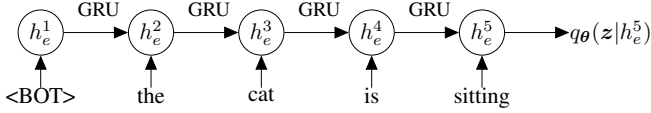


Figure 2. The encoder is implemented with GRU. Each word is represented by a parametric embedding. Given the input sequence, the encoder maps the last hidden state to the mean and variance of a Gaussian posterior over latent codes.

where  $p_\theta(x^k|\cdot)$  is a categorical distribution over  $|\mathcal{V}|$  possible tokens for the  $k^{th}$  element in  $\mathbf{x}$ . According to the model (Fig. 1), generating a sentence  $\mathbf{x}$  with latent code  $\mathbf{z}$  entails sampling each token from a distribution conditioned on the latent code and previously sampled tokens. Tokens are modelled with a parametric embedding. Conditioning the distribution over  $\mathbf{z}$  entails concatenating  $\mathbf{z}$  to the input embeddings per token (Bowman et al., 2015).

The encoder  $q_\theta(\mathbf{z}|\mathbf{x})$  is the posterior distribution over the latent variable  $\mathbf{z}$ , conditioned on a sequence  $\mathbf{x}$ . We take this to be Gaussian whose mean and diagonal covariance are specified by mappings  $\mu_\theta$  and  $\sigma_\theta$ :

$$q_\theta(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu_\theta(\mathbf{x}), \sigma_\theta(\mathbf{x})) \quad (3)$$

Linear mappings  $\mu_\theta$  and  $\sigma_\theta$  are computed from the last hidden state of a GRU (see Fig. 2).

## 2.2. MIM Learning Objective

The Mutual Information Machine (MIM) (Livne et al., 2019) is a versatile generative LVM which can be used for representation learning, and sample generation. MIM learns a model with high mutual information between observations and latent codes, and is robust against posterior collapse.

The MIM framework begins with two *anchor* distributions,  $\mathcal{P}(\mathbf{x})$  and  $\mathcal{P}(\mathbf{z})$ , for observations and the latent space, from which one can draw samples. They are fixed and not learned. MIM also has a parameterized encoder-decoder pair,  $q_\theta(\mathbf{z}|\mathbf{x})$  and  $p_\theta(\mathbf{x}|\mathbf{z})$ , and parametric marginal distributions  $q_\theta(\mathbf{x})$  and  $p_\theta(\mathbf{z})$ . These parametric elements define joint encoding and decoding *model* distributions:

$$q_\theta(\mathbf{x}, \mathbf{z}) = q_\theta(\mathbf{z}|\mathbf{x}) q_\theta(\mathbf{x}), \quad (4)$$

$$p_\theta(\mathbf{x}, \mathbf{z}) = p_\theta(\mathbf{x}|\mathbf{z}) p_\theta(\mathbf{z}). \quad (5)$$

MIM learning entails the minimization of the cross-entropy between a sample distribution and the model encoding and decoding distributions (Livne et al., 2019). This simple loss constitutes a variational upper bound on a regularized Jensen-Shannon divergence, resembling VAE in which a model distribution matches samples from a sample distribution via KL divergence minimization (Zhao et al., 2018). Fundamentally, MIM learning differs from VAE learning, with the former being an upper bound on the joint

cross-entropy, while the latter being an upper bound on the marginal cross-entropy of the observations.

MIM requires sampling from the decoder during training, which can be slow for sequential computational models. For language modeling we therefore use A-MIM learning, a MIM variant that minimizes a loss defined on the encoding and decoding distributions, with samples drawn from an encoding *sample* distribution, denoted  $\mathcal{M}_S^q(\mathbf{x}, \mathbf{z})$ ; i.e.,

$$\mathcal{M}_S^q(\mathbf{x}, \mathbf{z}) = q_\theta(\mathbf{z}|\mathbf{x}) \mathcal{P}(\mathbf{x}). \quad (6)$$

The A-MIM loss is defined as follows,

$$\begin{aligned} \mathcal{L}_{\text{A-MIM}}(\theta) &= \frac{1}{2} (CE(\mathcal{M}_S^q(\mathbf{x}, \mathbf{z}), q_\theta(\mathbf{x}, \mathbf{z})) \\ &\quad + CE(\mathcal{M}_S^q(\mathbf{x}, \mathbf{z}), p_\theta(\mathbf{x}, \mathbf{z}))) \\ &\geq H_{\mathcal{M}_S^q}(\mathbf{x}) + H_{\mathcal{M}_S^q}(\mathbf{z}) - I_{\mathcal{M}_S^q}(\mathbf{x}; \mathbf{z}), \end{aligned} \quad (7)$$

where  $CE(\cdot, \cdot)$  is cross-entropy,  $H_{\mathcal{M}_S^q}(\cdot)$  is information entropy over distribution  $\mathcal{M}_S^q$ , and  $I(\cdot; \cdot)$  is mutual information. Minimizing  $\mathcal{L}_{\text{A-MIM}}(\theta)$  learns a model with a consistent encoder-decoder, high mutual information, and low marginal entropy (Livne et al., 2019). The A-MIM loss is in fact a variational upper bound on the joint entropy of the encoding sample distribution.

## 2.3. Implicit and Explicit Model Marginals

To complete the model specification, we define the model marginals  $q_\theta(\mathbf{x})$  and  $p_\theta(\mathbf{z})$ . We call the marginal *explicit* when we can evaluate the probability of a sample under the corresponding distribution, and *implicit* otherwise.

Examples of explicit marginals are a Gaussian for  $p_\theta(\mathbf{z})$ , and an auto-regressive distribution for  $q_\theta(\mathbf{x})$ . They enable evaluation of the probability of a sample straightforwardly. However, the inductive bias in such distributions, or the architecture, can lead to a challenging optimization problem.

An implicit marginal distribution can be defined via marginalization of a joint distribution. To help encourage consistency, and avoid introducing more model parameters, one can define model marginals in terms of the sample distributions, like  $\mathcal{M}_S^q(\mathbf{x}, \mathbf{z})$  above (Bornschein et al., 2015; Livne et al., 2019; Tomczak & Welling, 2017). They allow one to share parameters between a marginal and the corresponding conditional distribution, and to reduce the inductive bias in the architecture. Unfortunately, evaluating the probability of a sample under an implicit distribution is intractable in general.

We define  $q_\theta(\mathbf{x})$  as a marginal over the decoder (Bornschein et al., 2015); i.e.,

$$q_\theta(\mathbf{x}) = \mathbb{E}_{\mathcal{P}(\mathbf{z})} [p_\theta(\mathbf{x}|\mathbf{z})], \quad (8)$$

where the latent anchor is defined to be a standard normal,  $\mathcal{P}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; 0, 1)$ . Similarly, the model density over latent

codes is defined as

$$p_{\theta}(\mathbf{z}) = \mathbb{E}_{\mathcal{P}(\mathbf{x})} [q_{\theta}(\mathbf{z}|\mathbf{x})] . \quad (9)$$

*i.e.*, the latent marginal is defined as the aggregated posterior, in the spirit of the VampPrior (Tomczak & Welling, 2017) and Exemplar VAE (Norouzi et al., 2020).

#### 2.4. Tractable Bounds to Loss

Given training data  $D = \{\mathbf{x}_i\}_{i=1}^N$ , the empirical loss is

$$\begin{aligned} \hat{\mathcal{L}}_{\text{A-MIM}}(\theta) = & -\frac{1}{2N} \sum_{\mathbf{x}_i} \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x}_i)} [\log q_{\theta}(\mathbf{z}|\mathbf{x}_i) q_{\theta}(\mathbf{x}_i)] \\ & - \frac{1}{2N} \sum_{\mathbf{x}_i} \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i|\mathbf{z}) p_{\theta}(\mathbf{z})] \end{aligned} \quad (10)$$

where  $\sum_{\mathbf{x}_i}$  denotes a sum over  $N$  fair samples drawn from  $\mathcal{P}(\mathbf{x})$ , as a MC approximation to expectation over  $\mathcal{P}(\mathbf{x})$ .

Unfortunately, the empirical loss in Eqn. (10) is intractable since we cannot evaluate the log-probability of the marginals  $p_{\theta}(\mathbf{z})$  and  $q_{\theta}(\mathbf{x})$ . In what follows we obtain a tractable empirical bound on the loss in Eqn. (10) for which, with one joint sample, we obtain an unbiased and low-variance estimate of the gradient using reparameterization (Kingma & Welling, 2013).

We first derive a tractable lower bound to  $\log q_{\theta}(\mathbf{x}_i)$ :

$$\begin{aligned} \log q_{\theta}(\mathbf{x}_i) & \stackrel{\text{(Eqn. 8)}}{=} \log \mathbb{E}_{\mathcal{P}(\mathbf{z})} [p_{\theta}(\mathbf{x}_i|\mathbf{z})] \\ & \stackrel{\text{(IS)}}{=} \log \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x}_i)} \left[ p_{\theta}(\mathbf{x}_i|\mathbf{z}) \frac{\mathcal{P}(\mathbf{z})}{q_{\theta}(\mathbf{z}|\mathbf{x}_i)} \right] \\ & \stackrel{\text{(JI)}}{\geq} \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x}_i)} \left[ \log \left( p_{\theta}(\mathbf{x}_i|\mathbf{z}) \frac{\mathcal{P}(\mathbf{z})}{q_{\theta}(\mathbf{z}|\mathbf{x}_i)} \right) \right] \end{aligned} \quad (11)$$

where the second and third lines are obtained using importance sampling and Jensen’s inequality. We remind the reader that  $q_{\theta}(\mathbf{x}_i)$  is a variational marginal that can depend on  $\mathbf{x}_i$ . Indeed, Eqn. (11) is the usual ELBO.

To derive a lower bound to  $\log p_{\theta}(\mathbf{z})$ , we begin with the following inequality,

$$\begin{aligned} \log \mathbb{E}_{\mathcal{P}(\mathbf{x})} [h(\mathbf{x}; \cdot)] & = \log \sum_i \mathcal{P}(\mathbf{x}_i) h(\mathbf{x}_i; \cdot) \\ & \geq \log \mathcal{P}(\mathbf{x}') h(\mathbf{x}'; \cdot) , \end{aligned} \quad (12)$$

for any sample  $\mathbf{x}'$ , any discrete distribution  $\mathcal{P}(\mathbf{x})$ , and any non-negative function  $h(\mathbf{x}; \cdot) \geq 0$ . The inequality in Eqn. (12) follows from  $\log a \geq \log b$  for  $a \geq b$ . Using this bound, we express a lower bound to  $p_{\theta}(\mathbf{z})$  as follows,

$$\begin{aligned} \log p_{\theta}(\mathbf{z}) & \stackrel{\text{(Eqn. 9)}}{=} \log \mathbb{E}_{\mathcal{P}(\mathbf{x})} [q_{\theta}(\mathbf{z}|\mathbf{x})] \\ & \stackrel{\text{(Eqn. 12)}}{\geq} \log q_{\theta}(\mathbf{z}|\mathbf{x}') + \log \mathcal{P}(\mathbf{x}') \end{aligned} \quad (13)$$

#### Algorithm 1 Learning parameters $\theta$ of sentenceMIM

```

1: while not converged do
2:    $D_{\text{enc}} \leftarrow \{\mathbf{x}_j, \mathbf{z}_j \sim q_{\theta}(\mathbf{z}|\mathbf{x})\mathcal{P}(\mathbf{x})\}_{j=1}^N$ 
3:    $\hat{\mathcal{L}}_{\text{MIM}}(\theta; D) = -\frac{1}{N} \sum_{i=1}^N ( \log p_{\theta}(\mathbf{x}_i|\mathbf{z}_i) + \frac{1}{2} (\log q_{\theta}(\mathbf{z}_i|\mathbf{x}_i) + \log \mathcal{P}(\mathbf{z}_i)) )$ 
4:    $\Delta\theta \propto -\nabla_{\theta} \hat{\mathcal{L}}_{\text{MIM}}(\theta; D)$  {Gradient computed through sampling using reparameterization}
5: end while
```

for any sample  $\mathbf{x}'$ . During training, given a joint sample  $\mathbf{x}_i, \mathbf{z}_i \sim q_{\theta}(\mathbf{z}|\mathbf{x})\mathcal{P}(\mathbf{x})$ , we choose  $\mathbf{x}' = \mathbf{x}_i$ .

Substituting Eqns. (11) and (13) into Eqn. (10) gives the final form of an upper bound on the empirical loss; *i.e.*,

$$\begin{aligned} \hat{\mathcal{L}}_{\text{A-MIM}} \leq & -\frac{1}{N} \sum_i \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x}_i)} [\log p_{\theta}(\mathbf{x}_i|\mathbf{z})] \\ & - \frac{1}{2N} \sum_i \mathbb{E}_{q_{\theta}(\mathbf{z}|\mathbf{x}_i)} [\log ( q_{\theta}(\mathbf{z}|\mathbf{x}_i)\mathcal{P}(\mathbf{z}) )] \\ & + \frac{1}{2} H_{\mathcal{P}}(\mathbf{x}) . \end{aligned} \quad (14)$$

We find an unbiased, low variance estimate of the gradient of Eqn. (14) with a single joint sample  $\mathbf{z}_i, \mathbf{x}_i \sim q_{\theta}(\mathbf{z}|\mathbf{x})\mathcal{P}(\mathbf{x})$  and reparameterization. The last term,  $H_{\mathcal{P}}(\mathbf{x})$ , is a constant, independent of model parameters and can therefore be ignored during optimization. The resulting learning process is described in Algorithm 1.

To better understand the proposed bounds, we note that A-MIM achieves good reconstruction by learning posteriors with relatively small variances (*i.e.*, relative to the distance between latent means). Our choice of  $\mathbf{x}' = \mathbf{x}_i$  exploits this, allowing good gradient estimation, and facilitating fast convergence. We further provide empirical evidence for these properties below in Fig. 3.

### 3. Experiments

#### 3.1. Datasets

(word level)	Sentences				
	Train	Valid.	Test	Vocab.	#words (avg.)
PTB	42068	3370	3761	9877	21 ± 10
Yahoo	100K	10K	10K	37165	76 ± 55
Yelp15	100K	10K	10K	19730	100 ± 51
WikiText-103	200K	10K	2185	89247	115 ± 60
Everything <sup>†</sup>	442067	33369	33760	105965	94 ± 60

Table 1. Dataset properties summary for Penn Tree Bank, Yahoo Answers and Yelp15, and sampled WikiText-103. Everything<sup>†</sup> is the union of all four datasets.

We show experimental results on four word level datasets described in Table 1, namely, Penn Tree Bank (Marcus et al.,



1993), Yahoo Answers and Yelp15 (following Yang et al. (2017)), and WikiText-103 (Merity et al., 2016). We use the Yahoo and Yelp15 datasets of Yang et al. (2017), which draw 100k samples for training, and 10k for validation and testing. For WT103 we draw 200k samples for training, 10k for validation, and retain the original test data. Empty lines and headers were filtered from the WT103 data.

### 3.2. Architecture and Optimization

Our auto-encoder architecture (Figs. 1 and 2) followed that proposed by Bowman et al. (2015). As is common, we concatenated  $z$  with the input to the decoder (*i.e.*, a "context", similar to He et al. (2019); Yang et al. (2017); Bowman et al. (2015)). We use the same architecture, parameterization, and latent dimensions for both sMIM and a VAE variant called sVAE, for comparison. We also trained deterministic auto-encoders with the same architecture, called sAE, by replacing the sampled latent code with the deterministic mean of the posterior (*i.e.*,  $z_i = \mathbb{E}_{z'} [q_\theta(z'|x_i)]$ ). Effectively, the only difference between these variants is the choice of loss function. Training times for all models are similar.

For PTB we trained models with 1 layer GRU, latent space dimensions of 16D, 128D, and 512D, a 512D hidden state, 300D word embeddings, and 50% embedding dropout. We trained all models with Adam (Kingma & Lei Ba, 2014) with initial learning rate  $lr = 10^{-3}$ . Training took less than 30 minutes on a single TITAN Xp 12G GPU. For Yahoo Answers, Yelp15, and WT103 we trained models with 1 layer GRU, latent space dimensions of 32D, 512D, 1024D, a 1024D hidden state, 512D word embeddings, and 50% embedding dropout. We trained these models with SGD (Sutskever et al., 2013), with initial  $lr = 5.0$ , and 0.25  $L_2$  gradient clipping. All model and optimization hyperparameters were taken from publicly available implementation of the method proposed by Bowman et al. (2015).

In all cases we use a learning rate scheduler that scaled the learning rate by 0.25 following two/one epochs (PTB/other datasets, respectively) with no improvement in the validation loss. We used a mini-batch size of 20 in all cases. Following (Sutskever et al., 2014) we feed the input in reverse to the encoder, such that the last hidden state in the encoder depends on the first word of the sentence.

We trained sVAEs with the regular ELBO, and with KL divergence annealing (denoted "+ kl"), where a scalar weight on the KL divergence term is increased from 0 to 1 over 10k mini-batches to lower the risk of posterior collapse (Bowman et al., 2015). We use no loss manipulation heuristics in the optimization of sMIM or sAE.

$z$ dim.	Enc. Recon. ↓	KL	Rand. Recon.	BLEU ↑	$ \theta $
sVAE (16)	105.24 (0.12)	1.6	105.91 (0.01)	0.124	11M
sVAE (128)	106.72 (0.12)	0.64	106.89 (0.01)	0.118	11M
sVAE (512)	108.52 (0.23)	0.41	108.88 (0.01)	0.116	12M
sAE (16)	91.86		163.9 (0.02)	0.348	11M
sAE (128)	67.56		113.61 (0.01)	0.589	11M
sAE (512)	62.08		102.93 (0.01)	0.673	12M
sMIM (16)	90.12 (0.03)		161.037 (0.02)	0.35	11M
sMIM (128)	67.35 (0.008)		136.2 (0.04)	0.61	11M
sMIM (512)	<b>59.23</b> (0.01)		133.74 (0.01)	<b>0.679</b>	12M
sMIM (1024) <sup>†</sup>	<b>26.43</b> (0.0)			<b>0.724</b>	179M

Table 2. Reconstruction results for **PTB** are averaged over 10 runs (stdev). Models<sup>†</sup> use extra training data. Reconstruction with a sample  $z \sim q_\theta(z|x)$  (Enc. Recon.) and with a random sample  $z \sim \mathcal{N}(z)$  (Rand. Recon.). An uninformative latent space will result in similar reconstruction values. see text for details.

$z$ dim.	Enc. Recon. ↓	KL	Rand. Recon.	BLEU ↑	$ \theta $
sVAE (32) + kl	401.63 (0.01)	31.86	425.92 (0.01)	0.274	40M
sVAE (512) + kl	379.93 (0.01)	4.19	385.76 (0.01)	0.18	43M
sVAE (1024) + kl	384.85 (0.01)	3.01	387.63 (0.01)	0.176	46M
sAE (32)	330.25		697.316 (0.0)	0.388	40M
sAE (512)	228.34		515.75 (0.0)	0.669	43M
sAE (1024)	222.7		503.87 (0.0)	<b>0.684</b>	46M
sMIM (32)	396.34 (0.0)		427.6 (0.0)	0.309	40M
sMIM (512)	220.03 (0.0)		600.29 (0.0)	0.673	43M
sMIM (1024)	<b>219.37</b> (0.0)		543.36 (0.0)	0.676	46M
sMIM (1024) <sup>†</sup>	<b>199.72</b> (0.0)			<b>0.686</b>	179M

Table 3. Reconstruction results for **Yelp15** are averaged over 10 runs. Models<sup>†</sup> use extra training data (See Table 3 for details).

### 3.3. Information Content in the Latent Code

Directly estimating the mutual information (MI) between the high-dimensional categorical observations  $x$  and the corresponding latent codes  $z$  is computationally expensive (Belghazi et al., 2018). Instead, we focus here on reconstruction, which is related to MI (Poole et al., 2019). We choose the reconstruction entropy (*Enc. Recon.*), which is the negative expected log-probability of the decoder, given a sample from the corresponding posterior. In addition, we show the reconstruction entropy when the latent code is sampled from a Gaussian prior (*Rand. Recon.*). The Gaussian has 0 mean and standard deviation fitted to the latent codes (see Table 11 in supplementary material). When the latent code conveys little information to the decoder, we expect Enc. Recon. and Rand. Recon. to be similar. When the decoder utilizes highly informative latent codes, we expect Rand. Recon. to be significantly larger (*i.e.*, worse). Finally we also show the 1-BLEU score, the fraction of words recovered in the sampled reconstruction.

Tables (2-4) show reconstruction results for PTB, Yelp15, Yahoo Answers. For all datasets but PTB, VAE learning with KL annealing was more effective than standard VAE

$z$ dim.	Enc. Recon. $\downarrow$	KL	Rand. Recon.	BLEU $\uparrow$	$ \theta $
sVAE (32) + kl	320.06 (0.04)	14.33	326.21 (0.01)	0.181	67M
sVAE (512) + kl	329.2 (0.06)	7.09	331.35 (0.01)	0.139	70M
sVAE (1024) + kl	334.41 (0.09)	5.52	335.83 (0.01)	0.131	73M
sAE (32)	293.75		487.45 (0.0)	0.372	67M
sAE (512)	222.34		375.38 (0.0)	0.624	70M
sAE (1024)	326.66		374.31 (0.0)	0.372	73M
sMIM (32)	290.37 (0.01)		555.82 (0.0)	0.387	67M
sMIM (512)	208.27 (0.0)		482.79 (0.01)	0.664	70M
sMIM (1024)	<b>205.81</b> (0.01)		475.16 (0.01)	<b>0.669</b>	73M
sMIM (1024) <sup>†</sup>	<b>178.82</b> (0.0)			<b>0.682</b>	179M

Table 4. Reconstruction results for **Yahoo Answers**, averaged over 10 runs. Models<sup>†</sup> use extra training data (See Table 3 for details).

$z$ dim.	Enc. Recon. $\downarrow$	KL	BLEU $\uparrow$	$ \theta $
sVAE (1024) + kl	481.66 (0.1)	12.65	0.165	153 M
sMIM (1024)	329.89 (0.02)		0.571	153 M
sMIM (1024) <sup>†</sup>	<b>313.66</b> (0.01)		<b>0.603</b>	179M

Table 5. Reconstruction results for **WT103** are averaged over 10 runs. Models<sup>†</sup> use extra training data. The superior reconstruction results of sMIM hold for longer sentences.

learning; due to the small size of PTB, annealing overfit. Model sMIM (1024)<sup>†</sup> is trained on all datasets (*i.e.*, PTB, Yahoo Answers, Yelp15 and WT103). The BLEU score is computed between test sentences and their reconstructed samples (higher is better), and  $|\theta|$  indicates the number of parameters in each model.

Tables (2-4) show that sMIM outperforms sVAE in reconstruction and BLEU score, and is comparable to sAE. In addition, the reconstruction of sMIM and sAE improves with more latent dimensions, showing that more information is captured by the latent codes, whereas sVAE shows the opposite trend due to posterior collapse (*i.e.*, encoder and random reconstructions are similar). Notice that sAE is more susceptible to over-fitting, as is evident in Table 4. WT103 results in Table 5 show that the superior reconstruction of sMIM also holds with longer sentences.

In summary, sMIM shows improved performance with more expressive architecture (higher latent dimension here), similar to sAE, while sVAE deteriorates due to posterior collapse. This suggests that sMIM could benefit from more powerful architectures like Transformers (Vaswani et al., 2017), without the need for posterior collapse-mitigating heuristics.

### 3.4. Posterior Collapse in VAE

The performance gap between sMIM and sVAE with identical architectures is due in part to posterior collapse in VAEs (*i.e.*, optimization is likely to have a role too), where the encoder has high posterior variance over latent codes, and hence low mutual information (cf. (Zhao et al., 2018; Alemi et al., 2017)); it coincides with the KL divergence term in the usual ELBO approaching zero in some or all

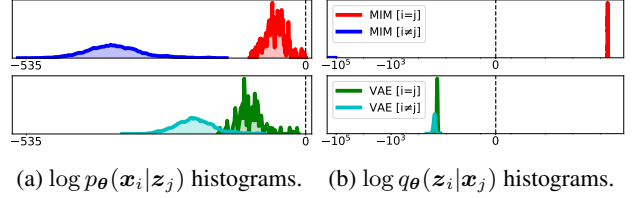


Figure 3. Histograms of log probabilities of test data for sMIM and sVAE trained on **PTB**: Overlap between curves indicates potential for poor reconstruction of input sentences. (a) Histograms of  $\log p_{\theta}(x_i | z_j)$  for  $z_j \sim q_{\theta}(z | x_j)$  when  $i = j$  (same input), and when  $i \neq j$  (when  $x_i$  is evaluated with the decoder distribution from a latent code associated with a different input sentence). (b) Histograms of  $\log q_{\theta}(z_i | x_j)$  for  $z_i \sim q_{\theta}(z | x_i)$ , when conditioned on the same input  $i = j$ , or a different input  $i \neq j$ .

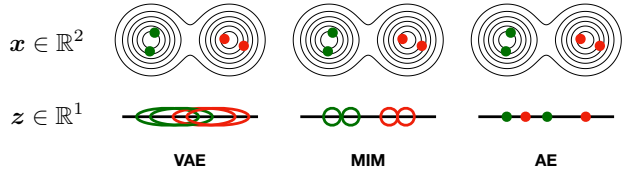


Figure 4. **Top**: level-sets of a 2D data distribution (*i.e.*, GMM with 2 modes). Red/green dots are samples. **Bottom**: the corresponding variance of the posterior, per sample, for MIM and VAE, and the corresponding mapping for AE. A semantically structured latent space will map latent samples from the same mode closer. A collapsed posterior (*i.e.*, VAE), might mix nearby modes, whereas MIM will have smaller posterior variance due to higher MI. A high entropy in the AE latent space might lead to non-semantic structure, where a perturbed red latent code might be reconstructed to the green mode, in contrast to MIM.

dimensions. In such cases, different sentences are mapped to similar regions of the latent space (see Fig. 3, bottom).

In contrast, given the high mutual information and reconstruction quality of sMIM, we only expect high decoding probability of a sentence when a latent code is sampled from the corresponding posterior. In other words, for sMIM, the posterior variances for different input sequences are relatively small compared to the distance between the posterior means (Fig. 3, top), allowing for accurate reconstruction.

The issue of posterior collapse becomes harder to mitigate as the dimension of the latent space increases, because the decoder becomes more expressive. As a consequence, language VAEs are typically limited to 32 dimensions or fewer (*e.g.*, He et al. (2019)), with only a few exceptions, such as Guu et al. (2017) which opted for 128 dimensions in a very particular problem settings. On the other hand, sMIM can easily scale up the latent dimension without issue.

### 3.5. Structure in the Latent Space

Here, we explore the structure in the learned representation. Table 6 shows the empirical entropy, estimated using NN

Dataset ( $z$ dim.)	sMIM	$\mathcal{N}$	sVAE	sAE
PTB (16D)	11.54 [ 0.67 ]	22.7	23.04 [ 1.02 ]	35.95 [ 1.18 ]
PTB (128D)	53.73 [ 0.55 ]	181.62	225.41 [ 1.24 ]	259.34 [ 1.26 ]
Yelp15 (32D)	32.22 [ 0.85 ]	45.4	50.01 [ 1.08 ]	73.03 [ 1.23 ]
Yelp15 (512D)	186.18 [ 0.56 ]	726.49	994.3 [ 1.37 ]	917.0 [ 1.3 ]
Yahoo (32D)	23.61 [ 0.73 ]	45.4	44.45 [ 1.03 ]	76.21 [ 1.26 ]
Yahoo (512D)	155.47 [ 0.48 ]	726.49	991.85 [ 1.37 ]	1003.26 [ 1.35 ]

Table 6. Empirical entropy of the latent codes, estimated with a NN entropy estimator. For comparison, column  $\mathcal{N}$  shows the entropy of a standard Normal in  $\mathbb{R}^d$  of corresponding latent dimension. In brackets is the ratio of the NN entropy to the entropy of an isotropic Gaussian fit to the latent codes. Ratios below 1 indicate that the latent codes are more clustered than a Gaussian, suggesting the low entropy for sMIM is not a simple consequence of down-scaling the latent codes.

entropy estimator Kraskov et al. (2004), of the latent codes for sMIM, sVAE, and sAE. Notice how the representation learned by sMIM has a significantly lower entropy. We note that sVAE is regularized to match a Gaussian, which is known to introduce smoother structure into the learned representation (see discussion by Bosc & Vincent (2020)).

With sMIM, we propose the use of information entropy minimization as an alternative regularization, which introduces meaningful structure without suffering from posterior collapse (see schematic plot in Fig. 4). Interestingly, neuroscientists have proposed that entropy minimization is an organizing principle in neural representations and information processing in the brain. Entropy minimization is viewed as allowing the agent to learn better to predict likely events (Friston, 2010; Barlow et al., 1972), compression/redundancy elimination (Barlow, 1961), and in terms of efficiency in energy consumption (Takagi, 2020).

By scaling the latent codes to reduce the variance of the aggregate posterior one can trivially reduce entropy with no benefit in terms of latent structure. To test whether this might be the case, we also fit an isotropic Gaussian to the latent codes (see Table 11 in supplementary materials), and show the ratio between the NN entropy and the fitted entropy in brackets. A ratio smaller than 1 suggests that the empirical entropy is more clustered than a Gaussian. Table 6 clearly shows that the lower empirical entropy cannot be explained by the scaling alone. We attribute the gap between the fitted and empirical entropies to clustering in the latent space, which MIM learning empirically demonstrates (Livne et al., 2019).

### 3.6. Reconstruction, Interpolation, and Perturbation

We further probe the structure in the learned representation, demonstrating that sMIM learns a dense, meaningful latent space. We present latent interpolation results in Table 7 for samples (*i.e.*, reviews) with the different ratings from Yelp5. Interpolation entails sampling  $x \sim p_\theta(x|z_\alpha)$  where  $z_\alpha$  is

#### 5 stars $\rightarrow$ 1 star

<BOT> awesome food , just awesome ! top notch beer selection . great staff . beer garden is great setting .
• awesome food , just top notch ! great beer selection . staff has great craft beer . top notch is that . <EOT>
• awesome food ! just kidding , beer selection is great . staff has trained knowledge on top . <EOT>
• cleanliness is awesome ! not only on their game , food . server was polite his hand sanitizer outside . <EOT>
• cleanliness is not on their patio . server was outside , kept running his hand sanitizer his hand . <EOT>
<BOT> cleanliness is not on their radar . outside patio was filthy , server kept running his hand thru his hair .

Table 7. Interpolation results between latent codes of input sentences (with gray) from **Yelp15** for sMIM (1024).

(D)	<BOT> the company did n't break out its fourth-quarter results
(M)	the company did n't break out its results <EOT>
(R)	the company did n't break out its fourth-quarter results <EOT>
(P)	the company did n't accurately out its results <EOT>

Table 8. Reconstruction results for sMIM (512) model trained on **PTB**. We denote: (D) Data sample; (M) Mean (latent) reconstruction; (R) Reconstruction; (P) Perturbed (latent) reconstruction.

interpolated at equispaced points between two randomly sampled latent codes,  $z_i \sim q_\theta(z|x_i)$ , and  $z_j \sim q_\theta(z|x_j)$ .

Next we show reconstruction, and perturbation results for for sMIM (512) trained on PTB. Table 8 shows four sentences: (D) the input sentence; (M) the mean reconstruction given the posterior mean  $z$ ; (R) a reconstruction given a random sample  $z$  from the posterior; and (P) a *perturbed reconstruction*, given a sample  $z$  from a Gaussian distribution with 10 times the posterior standard deviation. The high mutual information learned by sMIM leads to good reconstruction, as clear in (M) and (R). sMIM also exhibits good clustering in the latent space, shown here by the similarity of (R) and (P).

### 3.7. Question-Answering

So far we have discussed abstract aspects of representations learned by sMIM, such as high mutual information, and low marginal entropy. To demonstrate the benefits of representations learned by sMIM, we consider a downstream task in which sMIM is pre-trained on Yahoo Answers, then used for question-answering on YahooCQA (Tay et al., 2017b), with no fine-tuning. YahooCQA comprises questions and 3-5 answers, where the first answer is from Yahoo Answer, and the additional answers are wrong. Let  $Q_i$  denote the  $i^{th}$  question, and let  $\{A_i^k\}_{k=1}^{K_i}$  be the  $K_i$  corresponding answers, ordered such that  $A_i^k$  has rank  $k$ . To match the format of QA pairs in Yahoo Answers, we compose question-answer pair  $Q_i^k$  by concatenating  $Q_i$ , "?", and  $A_i^k$ .

For question-answering with sMIM we use the following procedure: For each question-answer we sample  $z_i^k \sim q_\theta(z|Q_i^k)$ , and a corresponding  $z_i^{unk} \sim q_\theta(z|Q_i^{unk})$  where

Model	P@1 ↑	MRR ↑
AP-CNN (dos Santos et al., 2016)	0.560	0.726
AP-BiLSTM (dos Santos et al., 2016)	0.568	0.731
HyperQA (Tay et al., 2017a)	<b>0.683</b>	0.801
sAE (32)	0.386	0.656
sAE (512)	0.579	0.814
sAE (1024)	0.519	0.767
sVAE (32) + kl *	0.531	0.776
sVAE (512) + kl *	0.494	0.747
sVAE (1024) + kl *	0.548	0.79
sMIM (32) ‡	0.558	0.737
sMIM (512) ‡	<b>0.683</b>	<b>0.818</b>
sMIM (1024) ‡	0.651	0.8
sAE (1024) †	0.574	0.812
sVAE (1024) †*	0.339	0.616
sMIM (1024) † ‡	<b>0.757</b>	<b>0.863</b>

Table 9. **YahooCQA** results for sMIM, AE, and single-task models. Results<sup>‡</sup> are averaged over 10 runs (stdev < 0.002). sMIM (1024)<sup>†</sup> is pre-trained on Everything dataset. sVAE\* results are based on the mean of the posterior, rather the sample (Bosc & Vincent, 2020). P@1 and MRR are defined in Sec. 3.7.

Q: <BOT> my brother is getting out on parole from navy jail where can i find a parole office in our area <UNK> , <UNK> ?
A: you can find out the county jail , or call your local police station . <EOT>
Q: <BOT> what continent has most deserts ?
A: the most notable is in the netherlands . <EOT>
Q: <BOT> how do u clear the history in the search field ?
A: u can find it in the search bar . <EOT>
Q: <BOT> what is the best question to ask ?
A: ask yourself ! <EOT>
Q: <BOT> need to find somewhere to sale baseball cards . ?
A: ebay <EOT>

Table 10. Sampled answers from **Yahoo Answers** sMIM (1024).

$Q_i^{unk}$  is simply  $Q_i$  concatenated with "?" and a sequence of <unk> tokens to represent the  $|A_i^k|$  unknown words of the answer. We then rank question-answer pairs according to the score  $S_i^k = ||z_i^{unk} - z_i^k|| / \sigma_i^{k,unk}$  where  $\sigma_i^{k,unk}$  is the standard deviation of  $q_\theta(z|Q_i^{unk})$ . In other words, we rank each question-answer pair according to the normalized distance between the latent code of the question with, and without, the answer. This score is similar to  $\log q_\theta(z_i^k|Q_i^{unk})$ , but without taking the log standard deviation into account.

As is common in the literature, Table 9 quantifies test performance using average precision ( $P@1 = \frac{1}{N} \sum_i \mathbb{1}(\text{rank}(A_i^1) = 1)$ ), and Mean Reciprocal Ranking ( $MRR = \frac{1}{N} \sum_i \frac{1}{\text{rank}(A_i^1)}$ ). As baselines, we consider best performing single-task models trained directly on YahooCQA (dos Santos et al., 2016; Tay et al., 2017a). Interestingly, sMIM (512), pre-trained on Yahoo Answers, exhibits state-of-the-art performance compared to these baselines. For an even larger sMIM model, pre-trained on all

of PTB, Yahoo Answers, Yelp15 and WT103, the question-answering performance of sMIM is even better (last row of Table 9).

The results for sVAE are based on the mean of the posterior rather than a random sample. This is a common heuristic in the NLP literature which has been proven useful for downstream tasks, but is problematic when considering the generative process which relies on the sample rather than the mean (Bosc & Vincent, 2020). Finally, as another point of comparison, we repeated the experiment with a deterministic sAE model (with  $\sigma_i^{k,unk} = 1$ ). In this case performance drops, especially average precision, indicating that the latent representations are not as meaningfully structured.

Importantly, sMIM can generate novel answers rather than simply ranking a given set of alternatives. To this end, we sample  $z_i^{unk} \sim q_\theta(z_i^k|Q_i^{unk})$ , as described above, followed by modified reconstruction  $\hat{Q}_i \sim p_\theta(x|z_i^{unk})$ . We modify the sampling procedure to be greedy (*i.e.*, top 1 token), and prevent the model from sampling the "<UNK>" token. We consider all words past the first "?" as the answer. (We also removed HTML tags (*e.g.*, "<br>").) Table 10 gives several selected answers. The examples were chosen to be short, and with appropriate (non-offensive) content. We note that we did not use any common techniques to manipulate the decoding distribution, such as beam search, Nucleus sampling (Holtzman et al., 2019), or sampling with temperature (Ackley et al., 1985). To the best of our knowledge, sMIM is the current state-of-the-art for a *single-task* model for YahooCQA, despite having simpler architecture and training procedure when compared to competing models.

## 4. Conclusions

This paper introduces sMIM, a new probabilistic auto-encoder for language modeling, trained with A-MIM learning. In particular, sMIM avoids posterior collapse, a challenging problem with VAEs applied to language, which enables the use of larger latent dimensions compared to sVAE, by orders of magnitude. While the reconstruction error is comparable to sAE, sMIM learns a latent representation with semantic structure, similar to sVAE, allowing for interpolation, perturbation, and sampling. In this sense, it achieves the best of both worlds: a semantically meaningful representation with high information content. We also use the structured latent representation for a downstream question-answering task on YahooCQA with state-of-the-art results. Importantly, the proposed framework has no hyperparameters in the loss, greatly simplifying the optimization procedure. In addition, sMIM benefits from a more expressive architecture, in contrast to sVAE, and demonstrates reduced susceptibility to over-fitting, compared to sAE. In future work, we will apply sMIM to more contemporary and powerful architectures like the Transformer.



## References

- Ackley, D. H., Hinton, G. E., and Sejnowski, T. J. A learning algorithm for boltzmann machines. *Cognitive Science*, 9 (1):147–169, 1985.
- Alemi, A. A., Poole, B., Fischer, I., Dillon, J. V., Saurous, R. A., and Murphy, K. An information-theoretic analysis of deep latent-variable models. *CoRR*, abs/1711.00464, 2017.
- Barlow, H. Possible principles underlying the transformations of sensory messages. In Rosenblith, W. (ed.), *Sensory Communication*. MIT Press, 1961.
- Barlow, H., Kaushal, T. P., and Mitchison, G. J. Finding minimum entropy codes. *Neural Computation*, 1:412–423, 1972.
- Belghazi, I., Rajeswar, S., Baratin, A., Hjelm, R. D., and Courville, A. MINE: Mutual information neural estimation. In *ICML*, 2018.
- Bornschein, J., Shabanian, S., Fischer, A., and Bengio, Y. Bidirectional Helmholtz machines. *CoRR*, abs/1506.03877, 2015.
- Bosc, T. and Vincent, P. Do sequence-to-sequence VAEs learn global features of sentences? In *EMNLP*, pp. 4296–4318. ACL, November 2020. doi: 10.18653/v1/2020.emnlp-main.350.
- Bowman, S. R., Vilnis, L., Vinyals, O., Dai, A. M., Józefowicz, R., and Bengio, S. Generating sentences from a continuous space. *CoRR*, abs/1511.06349, 2015.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. 2020.
- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1724–1734, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1179. URL <https://www.aclweb.org/anthology/D14-1179>.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. In *ACL*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://www.aclweb.org/anthology/N19-1423>.
- Dinh, L., Sohl-Dickstein, J., and Bengio, S. Density estimation using real nvp. *ICLR*, 2017.
- dos Santos, C. N., Tan, M., Xiang, B., and Zhou, B. Attentive pooling networks. *CoRR*, abs/1602.03609, 2016.
- Fang, L., Li, C., Gao, J., Dong, W., and Chen, C. Implicit deep latent variable models for text generation. In *EMNLP*, 2019.
- Friston, K. The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*, 11(2):127–138, 2010.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014.
- Guu, K., Hashimoto, T. B., Oren, Y., and Liang, P. Generating sentences by editing prototypes. *ACL*, 6:437–450, 2017.
- He, J., Spokoyny, D., Neubig, G., and Berg-Kirkpatrick, T. Lagging inference networks and posterior collapse in variational autoencoders. In *ICLR*, 2019.
- Hinton, G. E. and Zemel, R. Autoencoders, minimum description length and helmholtz free energy. In Cowan, J., Tesauro, G., and Alspector, J. (eds.), *NIPS*, volume 6, pp. 3–10. Morgan-Kaufmann, 1994.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, November 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL <http://dx.doi.org/10.1162/neco.1997.9.8.1735>.
- Holtzman, A., Buys, J., Forbes, M., and Choi, Y. The curious case of neural text degeneration. *CoRR*, abs/1904.09751, 2019.
- Kingma, D. P. and Lei Ba, J. ADAM: A method for stochastic optimization. In *ICLR*, 2014.
- Kingma, D. P. and Welling, M. Auto-Encoding Variational Bayes. In *ICLR*, 2013.
- Kraskov, A., Stögbauer, H., and Grassberger, P. Estimating mutual information. *Phys. Rev. E*, 69: 066138, Jun 2004. doi: 10.1103/PhysRevE.69.066138. URL <https://link.aps.org/doi/10.1103/PhysRevE.69.066138>.
- Kruengkrai, C. Better exploiting latent variables in text modeling. *ACL*, pp. 5527–5532, 2019.

- Li, C., Gao, X., Li, Y., Li, X., Peng, B., Zhang, Y., and Gao, J. Optimus: Organizing sentences via pre-trained modeling of a latent space. In *EMNLP*, 2020.
- Li, R., Li, X., Lin, C., Collinson, M., and Mao, R. A stable variational autoencoder for text modelling. In *INLG*, 2019a.
- Li, R., Li, X., Lin, C., Collinson, M., and Mao, R. A stable variational autoencoder for text modelling. *INLG*, pp. 594–599, 2019b.
- Livne, M., Swersky, K., and Fleet, D. J. MIM: Mutual Information Machine. *arXiv e-prints*, 2019.
- Marcus, M. P., Marcinkiewicz, M. A., and Santorini, B. Building a large annotated corpus of English: The Penn Treebank. *Comput. Linguist.*, 19(2):313–330, June 1993. ISSN 0891-2017. URL <http://dl.acm.org/citation.cfm?id=972470.972475>.
- Merity, S., Xiong, C., Bradbury, J., and Socher, R. Pointer sentinel mixture models. *CoRR*, abs/1609.07843, 2016. URL <http://arxiv.org/abs/1609.07843>.
- Merity, S., Keskar, N. S., and Socher, R. Regularizing and optimizing LSTM language models. *CoRR*, abs/1708.02182, 2017. URL <http://arxiv.org/abs/1708.02182>.
- Merity, S., Keskar, N. S., and Socher, R. Regularizing and optimizing LSTM language models. In *ICLR*, 2018.
- Norouzi, S., Fleet, D., and Norouzi, M. Exemplar VAE: Linking generative models, nearest neighbor retrieval, and data augmentation. In *NeurIPS*, 2020.
- Oord, A. v. d., Kalchbrenner, N., and Kavukcuoglu, K. Pixel recurrent neural networks. *International Conference on Machine Learning*, 2016.
- Poole, B., Ozair, S., Van Den Oord, A., Alemi, A., and Tucker, G. On variational bounds of mutual information. In *International Conference on Machine Learning*, pp. 5171–5180. PMLR, 2019.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. Language models are unsupervised multitask learners. 2019.
- Rae, J. W., Dyer, C., Dayan, P., and Lillicrap, T. P. Fast parametric learning with activation memorization. *CoRR*, abs/1803.10049, 2018. URL <http://arxiv.org/abs/1803.10049>.
- Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference in deep tve Models. In *ICML*, 2014.
- Sutskever, I., Martens, J., Dahl, G., and Hinton, G. On the importance of initialization and momentum in deep learning. In Dasgupta, S. and McAllester, D. (eds.), *ICML*, volume 28 of *Proceedings of Machine Learning Research*, pp. 1139–1147, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.
- Sutskever, I., Vinyals, O., and Le, Q. V. Sequence to sequence learning with neural networks. In *NIPS*, NIPS, pp. 3104–3112, Cambridge, MA, USA, 2014. MIT Press. URL <http://dl.acm.org/citation.cfm?id=2969033.2969173>.
- Takagi, K. Principles of mutual information maximization and energy minimization affect the activation patterns of large scale networks in the brain. *Frontiers in Computational Neuroscience*, 13:86, 2020. ISSN 1662-5188.
- Tay, Y., Luu, A. T., and Hui, S. C. Enabling efficient question answer retrieval via hyperbolic neural networks. *CoRR*, abs/1707.07847, 2017a. URL <http://arxiv.org/abs/1707.07847>.
- Tay, Y., Phan, M. C., Luu, A. T., and Hui, S. C. Learning to rank question answer pairs with holographic dual LSTM architecture. In *SIGIR*, pp. 695–704, 2017b. doi: 10.1145/3077136.3080790. URL <http://doi.acm.org/10.1145/3077136.3080790>.
- Tomczak, J. M. and Welling, M. VAE with a vampprior. *CoRR*, abs/1705.07120, 2017. URL <http://arxiv.org/abs/1705.07120>.
- Vahdat, A. and Kautz, J. NVAE: A deep hierarchical variational autoencoder. In *NeurIPS*, 2020.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. In *NIPS*, pp. 5998–6008, 2017.
- Wang, D., Gong, C., and Liu, Q. Improving neural language modeling via adversarial training. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *ICML*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6555–6565, Long Beach, California, USA, 09–15 Jun 2019. PMLR. URL <http://proceedings.mlr.press/v97/wang19f.html>.
- Yang, Z., Hu, Z., Salakhutdinov, R., and Berg-Kirkpatrick, T. Improved variational autoencoders for text modeling using dilated convolutions. *CoRR*, abs/1702.08139, 2017. URL <http://arxiv.org/abs/1702.08139>.
- Zhao, S., Song, J., and Ermon, S. A Lagrangian perspective on latent variable generative models. *UAI*, Jul 2018.

## A. Distribution of Sentence Lengths

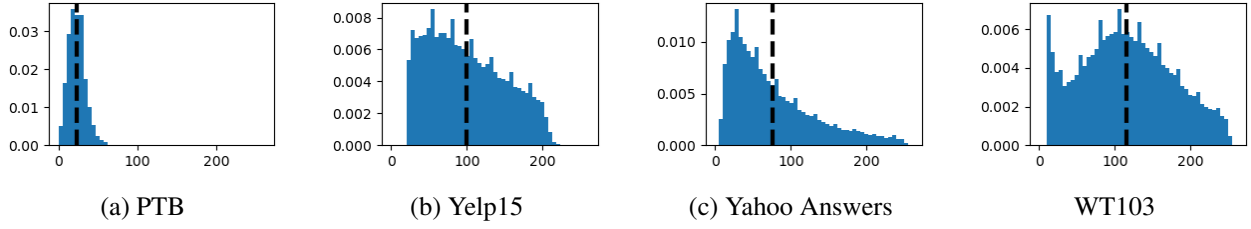


Figure 5. Here we present histograms of sentence lengths per dataset. The dashed line is the average sentence length.

Fig. 5 shows histograms of sentence lengths. Notice that PTB sentences are significantly shorter than other datasets. As a result, sMIM is somewhat better able to learn a representation that is well suited for reconstruction. Other datasets, with longer sentences, are more challenging, especially with the simple architecture used here (*i.e.*, 1 later GRU). We believe that implementing sMIM with an architecture that better handles long-term dependencies (*e.g.*, transformers) might help.

## B. Comparison of Reconstruction in MIM and VAE

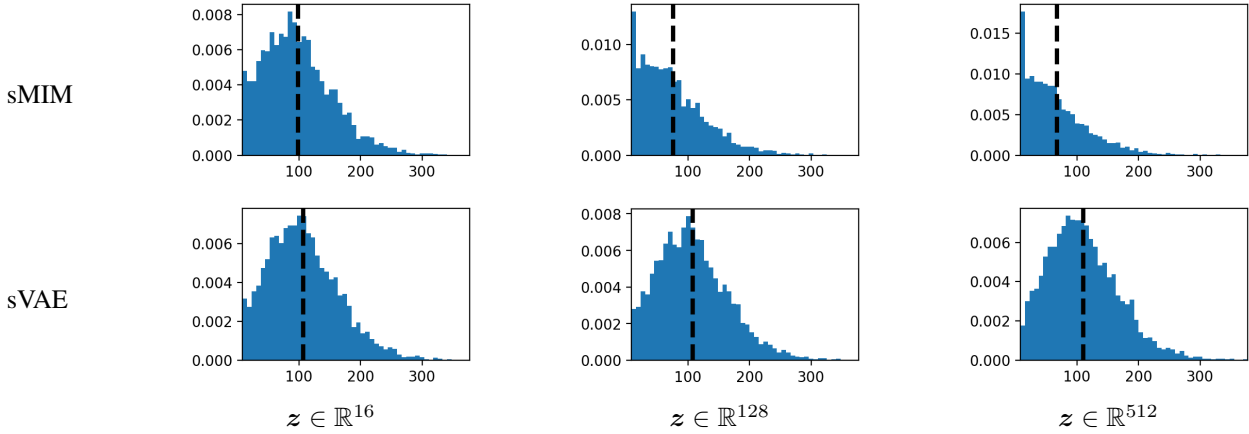


Figure 6. Histograms of reconstruction for sMIM and sVAE versus latent dimension for **PTB**. Dashed black line is the mean.

Figures 6-8 depict histograms of reconstruction values for sentences, for sVAE and sMIM with different latent dimensions. While a less expressive sMIM behaves much like sVAE, the difference is clearer as the expressiveness of the model increases. Here, sVAE does not appear to effectively use the increased expressiveness for better modelling. We hypothesize that the added sVAE expressiveness is used to better match the posterior to the prior, resulting in posterior collapse. sMIM uses the increased expressiveness to increase mutual information.

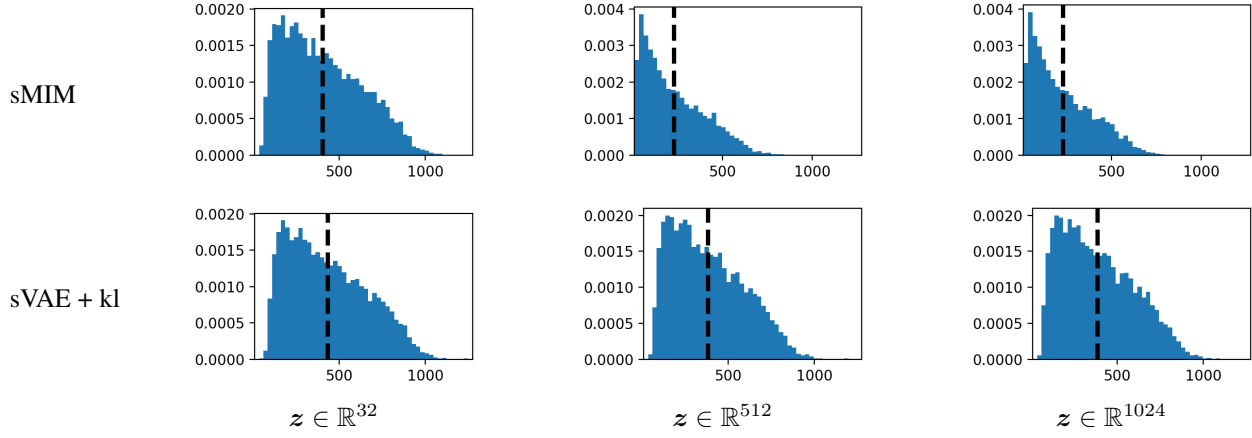


Figure 7. Histograms of reconstruction for sMIM and sVAE versus latent dimension for **Yelp15**. Dashed black line is the mean.

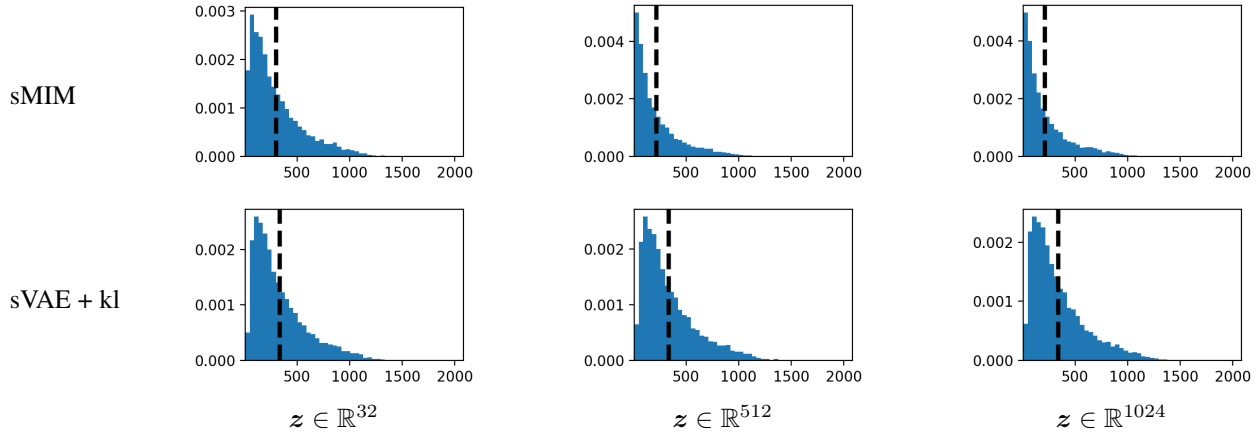


Figure 8. Histograms of reconstruction for sMIM and sVAE versus latent dimension for **Yahoo Answers**.

### C. Empirical Latent Entropy

Table 11 provides the entropy of an isotropic Gaussian that is fitted to the latent codes, and the standard deviation of the fitted Gaussian [entropy / stdev]. The values are used to compute the ratio presented in the main paper.



## SentenceMIM

Dataset ( $z$ dim.)	sMIM	$\mathcal{N}$	sVAE	sAE
PTB (16D)	[ 17.22 / 0.5 ]	22.7	[ 22.51 / 0.97 ]	[ 30.49 / 2.64 ]
PTB (128D)	[ 97.39 / 0.26 ]	181.62	[ 181.29 / 0.99 ]	[ 206.22 / 1.46 ]
Yelp15 (32D)	[ 38.06 / 0.63 ]	45.4	[ 46.31 / 1.05 ]	[ 59.16 / 0.99 ]
Yelp15 (512D)	[ 333.45 / 0.21 ]	726.49	[ 726.15 / 0.99 ]	[ 705.14 / 0.91 ]
Yahoo (32D)	[ 32.13 / 0.43 ]	45.4	[ 43.3 / 0.87 ]	[ 60.45 / 2.56 ]
Yahoo (512D)	[ 326.17 / 0.2 ]	726.49	[ 724.75 / 0.99 ]	[ 744.24 / 1.07 ]

Table 11. In brackets is the entropy of an isotropic Gaussian fitted to the latent codes, and the corresponding average standard deviation [ entropy / stdev ]. For comparison, column  $\mathcal{N}$  shows the entropy of a standard Normal in  $\mathbb{R}^d$  of a corresponding latent dimension. Our goal here is to rule out simple down-scaling as the cause for the low entropy in sMIM.

## D. Additional Results

### D.1. Reconstruction

	sMIM (512)	sMIM (1024) <sup>†</sup>
(D)	<BOT> there was no panic	
(M)	there was no panic <EOT>	there was no panic <EOT>
(R)	there was no orders <EOT>	there was no panic <EOT>
(P)	there was no panic <EOT>	there was no shortage panic <EOT>
(AE)	there was no panic <EOT>	
(D)	<BOT> the company did n't break out its fourth-quarter results	
(M)	the company did n't break out its fourth-quarter results <EOT>	the company did n't break out its results results <EOT>
(R)	the company did n't break out its results <EOT>	the company did n't break out its results <EOT>
(P)	the company did n't break out its fourth-quarter results <EOT>	the company did n't break out its results results <EOT>
(AE)	the company did n't break out results <EOT>	
(D)	<BOT> it had planned a strike vote for next sunday but that has been pushed back indefinitely	
(M)	it had a weakening for promotional planned but that has pushed aside back but so far away <EOT>	it had planned planned a planned for next week but that continues has been pushed back pushed <EOT>
(R)	it had a planned strike for energy gifts but so that has planned airlines but block after six months <EOT>	it had planned a strike planned for next sunday but that has been pushed back culmination pushed <EOT>
(P)	it had a strike with stateswest airlines but so that it has slashed its spending but so far said he would be subject by far <EOT>	it had planned a strike for hardcore but has been pushed every year that leaves back <EOT>
(AE)	it had been a five-year vote but for a week that drilling humana strike back back has planned back <EOT>	

Table 12. Reconstruction results for models trained on **PTB**. We denote: (D) Data sample; (M) Mean (latent) reconstruction; (R) Reconstruction; (P) Perturbed (latent) reconstruction; (AE) Reconstruction of AE.

Here we provide reconstruction results for PTB (Fig. 12), Yelp15 (Fig. 13), and Yahoo Answers (Fig. 14). Each figure shows (D) Data sample; (M) Mean (latent) reconstruction (*i.e.*,  $z_i = \mathbb{E}[q_\theta(z|x_i)]$ ); (R) Reconstruction (*i.e.*,  $z_i \sim q_\theta(z|x_i)$ ); (P) Perturbed (latent) reconstruction (*i.e.*,  $z_i \sim q_\theta(z|x_i; \mu_i, 10\sigma_i)$ ); (AE) Reconstruction of AE. We compare the best performing sMIM model to an AE with the same architecture, and to sMIM (1024) <sup>†</sup> (*i.e.*, the model trained on the Everything dataset).

Interestingly, AEs tend to perform worse for longer sentences, when compared to sMIM. We attribute this to the higher latent entropy, which leads to non-semantic errors (*i.e.*, nearby latent codes are less similar compared to MIM). Another interesting point is how the reconstruction (R), is better in many cases than the reconstruction given the mean latent code from the encoder (M) (*i.e.*, which have the highest probability density). We attribute that to the fact that most probability mass in a high dimensional Gaussian in  $d \gg 1$  dimensional space and  $\sigma$  standard deviation is concentrated in around a sphere of radius  $r \approx \sigma\sqrt{d}$ . As a result the probability mass around the mean is low, and sampling from the mean is less likely to represent the input sentence  $x_i$ . This also explains how perturbations of up to 10 standard deviations might result in good reconstructions. Finally, we point how sMIM (1024) <sup>†</sup>, trained on Everything, does a better job handling longer sentences.

# SentenceMIM

	sMIM (1024)	sMIM (1024) <sup>†</sup>
(D)	<b>(3 stars)</b> <BOT> decent price . fast . ok staff ... but it is fast food so i ca n't rate any higher than 3 .	
(M)	decent italians . fast . price ok ... but it is higher than any other fast food i ca n't rate so higher rate jusqu . <EOT>	decent oxtail . ok . fast price ... but staff it is so fast i ca n't rate any food 3 . <EOT>
(R)	decent price . superior . decent staff ... but ok fast food is n't so it i ' d rate higher any higher quality than 3 . <EOT>	decent price . fast staff . fast ok ... but it is so fast food i rate 3 higher than any . <EOT>
(P)	decent price . ok . fast food ... but it is ok . so i ca n't rate any higher rate as fast food is marginal . <EOT>	decent price . fast . wu ... fast food ! but it staff so ok i ca n't rate 3 stars . . <EOT>
(AE)	decent price . fast staff . ok ... but it is fast food so i ca n't rate any rate than 3 . <EOT>	
(D)	<b>(4 stars)</b> <BOT> excellent wings . great service . 100 % smoked wings . great flavor . big meaty . i will definitely be back . okra is great too .	
(M)	excellent wings . great service . 100 % wings . big meaty wings . great flavor . i definitely will be back . lake is great too . <EOT>	excellent service . great wings . 100 % superior . great flavor . great fries . definitely will be back . i had too big fat . <EOT>
(R)	excellent wings . great service . 100 % wings . wings flavor . definitely great . 100 % . i will be back . <EOT>	excellent service . great flavor . 100 % wings . excellent . great big guts . definitely will be back from . i had great wings . <EOT>
(P)	excellent wings . great service . wings flavours wings . 100 % big . mmmmm overwhelmed . i ' m definitely hooked . bye disgusted is great but will be back . i definitely go . <EOT>	great burger . excellent service . 100 % fat bowls . great carnitas . great flavor . i will definitely be back . i avoid too late . <EOT>
(AE)	excellent excellent . great service . 100 % wings . 100 % big burritos . 100 % . i will definitely be back . great too too is ultra <EOT>	
(D)	<b>(5 stars)</b> <BOT> delicious ! the sandwiches are really good and the meat is top quality . it ' s also nice grabbing an exotic item from the shelf for dessert .	
(M)	delicious ! the meat really are good and the quality is nice . it ' s also tempting top notch lovers from the roasters an item top . <EOT>	delicious ! the sandwiches are really good and the quality is top notch . it ' s an exotic item popping also generates from the top spices . <EOT>
(R)	delicious ! the sandwiches are really good and the meat is quality . it ' s also nice dessert for shipping from the top floor an unhygienic machine . <EOT>	delicious ! the sandwiches are really good and the quality is top notch . it ' s also charging an item assortment from the grocery store for dessert . <EOT>
(P)	delicious sandwiches ! the servers are really good and the quality is top notch . it ' s also an item for meat quality memories . <EOT>	who ! the meat are really good and the quality is top notch ' s . it also seems top notch item has yet and an unexpected range for the pistachio . i do cross like john tomatoes from my experience . <EOT>
(AE)	delicious ! the sandwiches are really good and the quality is top notch . it ' s also caught meat also fixing an item from the top for nice hash . <EOT>	

Table 13. Reconstruction results for models trained on **Yelp15**. We denote: (D) Data sample; (M) Mean (latent) reconstruction; (R) Reconstruction; (P) Perturbed (latent) reconstruction; (AE) Reconstruction of AE.

	sMIM (1024)	sMIM (1024) <sup>†</sup>
(D)	<b>(Sports)</b> <BOT> are you regular or goofy ? regularly goofy	
(M)	are you regular or regular ? regular <EOT>	are you regular or regularly ? regular johnny <EOT>
(R)	are you regular regular or nintendo ? regular icecream <EOT>	are you regular or regularly ? regularly gethsemane <EOT>
(P)	are you or regular worms regular ? regular goldfish by benjamin <EOT>	are you regular or early regularly regularly regularly <EOT>
(AE)	are you sex or two frustrated <EOT>	
(D)	<b>(Health)</b> <BOT> how do you start to like yourself ? i was taught by my parents .	
(M)	how do you start to like yourself ? i would like to meet my parents by . <EOT>	how do you start to like yourself ? i was taught by my parents . <EOT>
(R)	how do you start to yourself like ? i was taught my parents by parents . <EOT>	how do you start to like yourself ? i was taught by my parents . <EOT>
(P)	how do you start to like yourself ? i am 27 by my self . <EOT>	how do you start to like yourself ? start by i was taught my foot . <EOT>
(AE)	how do you like to after by christmas day ? i like to aid my boss by my brother and state ! <EOT>	
(D)	<b>(Business &amp; Finance)</b> <BOT> how can i find someone in spain ? i'm in spain today , what do you want ?	
(M)	how can i find someone in spain ? i'm in harlem limo , now what do you want ? <EOT>	how can i find someone in spain ? spain in spain ? i'm talking , what did you want ? <EOT>
(R)	where can i find someone in spain ? in spain today , what do you want ? <EOT>	how can i find someone in spain ? spain in spain today , what do you want ? <EOT>
(P)	how can i find someone in stone ? in nassau i'm sure civilian , what ? you want today ! <EOT>	how can i find someone in spain ? i'm in spain today ? what maytag , do you think ? <EOT>
(AE)	how can i find someone in africa investment , ca ? working 6.0 in future with susan toughie <EOT>	

Table 14. Reconstruction results for models trained on **Yahoo Answers**. We denote: (D) Data sample; (M) Mean (latent) reconstruction; (R) Reconstruction; (P) Perturbed (latent) reconstruction; (AE) Reconstruction of AE.

## D.2. Interpolation

sMIM (512)	sMIM (1024) <sup>†</sup>
<BOT> thanks to modern medicine more couples are growing old together	
<ul style="list-style-type: none"> <li>• to growing small businesses are growing more rapidly growing &lt;EOT&gt;</li> <li>• growing to more areas are growing preventing black trends &lt;EOT&gt;</li> <li>• growing to the growing industry are growing more rapidly growing than &lt;EOT&gt;</li> <li>• growing to the exact industry has been growing more sophisticated six months &lt;EOT&gt;</li> <li>• politics the growing issue are not to mention closely although other prospective products &lt;EOT&gt;</li> <li>• the system is growing enough to make not radical an article &lt;EOT&gt;</li> <li>• the system is reducing compliance not to consider an article &lt;EOT&gt;</li> <li>• the system is the problem system not an effective &lt;EOT&gt;</li> <li>• the system is the system not knowing an individual &lt;EOT&gt;</li> <li>• the system is the system not an encouraging problem &lt;EOT&gt;</li> </ul>	<ul style="list-style-type: none"> <li>• thanks to modern medicine more modern couples are growing together than &lt;EOT&gt;</li> <li>• thanks to modern cancer more are growing peaceful couples form &lt;EOT&gt;</li> <li>• thanks to medicine rosen modern more are growing together governing &lt;EOT&gt;</li> <li>• thanks to moolah the modern premises are more sensitive together &lt;EOT&gt;</li> <li>• programm thanks to the cutbacks schedules is not an church system &lt;EOT&gt;</li> <li>• humana remains the loyalty to instituting dynamic is an orthodox montage &lt;EOT&gt;</li> <li>• the strategies is not paying the non-food system an individual member &lt;EOT&gt;</li> <li>• the system is not the individual problem member an can &lt;EOT&gt;</li> <li>• the system is not the individual problem an individual member &lt;EOT&gt;</li> <li>• the system is not the individual problem an individual member &lt;EOT&gt;</li> </ul>
<BOT> the system is the problem not an individual member	
<ul style="list-style-type: none"> <li>• the system is the system not an investment fund &lt;EOT&gt;</li> <li>• the system is the problem not an office &lt;EOT&gt;</li> <li>• the system is not the problem for an individual &lt;EOT&gt;</li> <li>• the system is not clear the veto &lt;EOT&gt;</li> <li>• the system is not encouraging to the securities &lt;EOT&gt;</li> <li>• xtra the system is not even critical &lt;EOT&gt;</li> <li>• sony denies the declines to secure &lt;EOT&gt;</li> <li>• everyone brought the stock to comment &lt;EOT&gt;</li> <li>• sony which declines to comment &lt;EOT&gt;</li> <li>• kellogg declines to induce itself &lt;EOT&gt;</li> </ul>	<ul style="list-style-type: none"> <li>• the system is the ringers not an individual member &lt;EOT&gt;</li> <li>• the system is not the problem an individual member &lt;EOT&gt;</li> <li>• the problem is not the indies system an individual &lt;EOT&gt;</li> <li>• the merksamer is not the problem system an individual &lt;EOT&gt;</li> <li>• mr . the herald is not an individual problem &lt;EOT&gt;</li> <li>• qintex producers is the president's to comment &lt;EOT&gt;</li> <li>• sony preferences itself is the bidding to comment &lt;EOT&gt;</li> <li>• sony sony itself is to comment &lt;EOT&gt;</li> <li>• sony sony itself to comment &lt;EOT&gt;</li> <li>• sony declines itself to sony &lt;EOT&gt;</li> </ul>
<BOT> sony itself declines to comment	

Table 15. Interpolation results between latent codes of input sentences (with gray) from **PTB**.

Here we provide interpolation results for PTB (Fig. 15), Yelp15 (Fig. 16), and Yahoo Answers (Fig. 17). We compare the best performing sMIM model to sMIM (1024) <sup>†</sup>. Interestingly, both models appear to have learned a dense latent space, with sMIM (1024) <sup>†</sup> roughly staying within the domain of each dataset. This is surprising since the latent space of sMIM (1024) <sup>†</sup> jointly represents all datasets.

sMIM (1024)	sMIM (1024) <sup>†</sup>
<b>(3 star)</b> <BOT> as bbq in phoenix goes - this is one of the better ones . get there early - they fill up fast !	
<ul style="list-style-type: none"> <li>• as in china phoenix - this is one of the better ones fast get . fill there early - they fill up early ! &lt;EOT&gt;</li> <li>• as far in san jose - this is one of the better ones . fast get up early ! there they fill up fast for u ! &lt;EOT&gt;</li> <li>• as pei wei goes in this phoenix - - one of the best ones . get there early ! they picked up fast food items is better . &lt;EOT&gt;</li> <li>• oxtail yo buffet in pittsburgh as the owners goes - better . this is not one of those fast food places . fill up there get the hot ! &lt;EOT&gt;</li> <li>• ah circle k ! not as bad in the food . thankfully - this one is one of the best bbq joints here ! service was fast friendly . &lt;EOT&gt;</li> <li>• ehh = ciders as the food goes . not bad for service ! - in many fast the only ones available is this . you can get better steak anywhere else ! &lt;EOT&gt;</li> <li>• bin spaetzle food not the best . wicked spoon ! service is brutal only fast for the hot mexican in lv . everything else on this planet as can you get . &lt;EOT&gt;</li> <li>• frankie food not soo the best . service = horrible ! only drawback frozen for these hike . everything you can pass on the juke planet . &lt;EOT&gt;</li> <li>• food not the best service . knocking only 99 cents ! for the hot buffet everything . beef &amp; broccoli on the vip polo you can pass . &lt;EOT&gt;</li> <li>• food not the best . service = horrible ! only plopped for the paella everything &amp; rum . you can find everything on the strip . &lt;EOT&gt;</li> </ul>	<ul style="list-style-type: none"> <li>• as in phoenix goes this is - better than one of the newest ones . get there early - they fill up fast ! &lt;EOT&gt;</li> <li>• as shore goes in phoenix - this is one of the better bbq . fast ! they get up there early - men dinner . &lt;EOT&gt;</li> <li>• as dean goes in phoenix this is the list of bbq . - one not goes fast - get there early ! they fill up fast . &lt;EOT&gt;</li> <li>• veal as rocks as this goes in the phoenix area . - one of food is not better quick enough they get . 2 enchiladas up ! &lt;EOT&gt;</li> <li>• kohrs as molasses as comparing goes in the food . not sure is one of this better ones - the only ones for fat . thumbs squeeze there ! &lt;EOT&gt;</li> <li>• omg = rainbow not as the food goes . congrats service ! this is one of the hot spots for only frozen hot - you can . eat on carts there . &lt;EOT&gt;</li> <li>• = frozen food ! not the best . only frozen hot as for you shall pick the ice cream - . loved everything else on wednesday ! &lt;EOT&gt;</li> <li>• = food not only the best . frozen service ! everything else for the frozen yogurt company . absolute hot tea during normal on as they can . &lt;EOT&gt;</li> <li>• = food not . the best frozen service ! only five stars for the water suppose . hot things you can smell on budget . &lt;EOT&gt;</li> <li>• food = not the best . frozen service ! only \$ 21 for the frozen hot chocolate . everything else can you tell on romance . &lt;EOT&gt;</li> </ul>
<b>(2 star)</b> <BOT> food = not the best . service = horrible ! only known for the frozen hot chocolate . everything else you can pass on .	
<ul style="list-style-type: none"> <li>• food not the best . fuck service only ! ! horrible cannolis for the fajitas unusual known . everything you can pass on graduate . &lt;EOT&gt;</li> <li>• food not suck . the best service ever ! just horrible everything for the frozen hot chocolate . you can probably survive on everything else . &lt;EOT&gt;</li> <li>• food = not ! service = the best . only organizations thing for chocolate lovers treats and green beans . everything you can taste on the planet . &lt;EOT&gt;</li> <li>• blech food ! not the best dish anywhere else . service = &lt;unk&gt; for the frozen hot chocolate and dessert bartenders ! everything you can only expect better at this shuffle . &lt;EOT&gt;</li> <li>• 32 words ! not amazing food . the best &lt;unk&gt; music and service they had can earned a better meal at xs . everything else on bill for me . &lt;EOT&gt;</li> <li>• snottsdale act ! ! rio mia &lt;unk&gt; at the food and wished you not a fan . delicious lunch &amp; dessert better choices for dessert but they had blackjack . &lt;EOT&gt;</li> <li>• husbands cher ! wish they had &lt;unk&gt; dessert at the bellagio and not a great lunch selection . food better tasting wise but sadly serves and dessert selection . &lt;EOT&gt;</li> <li>• soooo ! pretzel panera &lt;unk&gt; they had at a better selection and the food sucked but nothing memorable a dessert . surely great value and better mayonnaise desserts . &lt;EOT&gt;</li> <li>• yummy ! wish they had &lt;unk&gt; at lunch and a dessert selection but a better value and great value than beef suggestion company . &lt;EOT&gt;</li> <li>• yummy ! wish they had &lt;unk&gt; dessert at lunch and a selection but a tiramisu better value and freshness value food taste better than ihop . &lt;EOT&gt;</li> </ul>	<ul style="list-style-type: none"> <li>• food = not the best . frozen hot service ! only website for the frozen hot chocolate . you can grab everything else on . &lt;EOT&gt;</li> <li>• food = not the best . frozen service ! only for five stars during the san francisco frozen chicken . everything else on could not give thumbs . &lt;EOT&gt;</li> <li>• food = not ! the frozen yogurt . service only best for you ate here twice although the frozen yogurt . delicious atmosphere on everything else . &lt;EOT&gt;</li> <li>• gelato food ! not sure the best . frozen seared only wish you can mix for the frozen hot chocolate frozen . service on and everything else explains . &lt;EOT&gt;</li> <li>• hilariously = ! food is not the best meal . hibachi cover service and they only wished a frozen yogurt for hot girl . better luck at &lt;unk&gt; and on the latter experience . &lt;EOT&gt;</li> <li>• blended ! wifey better food ! the service is not frozen hot . they redeemed a &lt;unk&gt; and only frozen someplace at horse's for frozen worms . &lt;EOT&gt;</li> <li>• wish ! methinks buffet is ingrediants at the &lt;unk&gt; food and a better tasting . they woulda frozen lunch but not memorable and satisfying tasting better ambiance . &lt;EOT&gt;</li> <li>• yummy ! wish they had &lt;unk&gt; at a buffet and netherlandish better tasting food . a renovation treasure and great value but not better than calories tasting . &lt;EOT&gt;</li> <li>• wish ! wish they had &lt;unk&gt; at 10am and a dessert selection but better food a better and better tasting selection . great value ! &lt;EOT&gt;</li> <li>• wish ! wish they had lunch at &lt;unk&gt; and a dessert fountain but better than a selection and great tasting food servings better tasting . &lt;EOT&gt;</li> </ul>
<b>(4 star)</b> <BOT> yummy ! wish they had <unk> at lunch and a better dessert selection but a great value and better tasting food than wicked spoon .	

Table 16. Interpolation results between latent codes of input sentences (with gray) from Yelp15.



sMIM (1024)	sMIM (1024) <sup>†</sup>
<b>(Business &amp; Finance)</b> <BOT> are u shy or outgoing ? both , actually	
<ul style="list-style-type: none"> <li>• are u or wishing vidio ? both , actually &lt;EOT&gt;</li> <li>• are u or stressed caffiene ? both , actually make a smile &lt;EOT&gt;</li> <li>• witch are u or how lucky ? both &lt;EOT&gt;</li> <li>• are u kidding or spraying ? both &lt;EOT&gt;</li> <li>• how does wile or are you ? to both use , instead like it . &lt;EOT&gt;</li> <li>• how do u choose to start or ? like i cant think , are actually better by my work . &lt;EOT&gt;</li> <li>• how do you start to alienate yourself ? i are like or drone , my actually feels . &lt;EOT&gt;</li> <li>• how do you start to yourself or like ? i like my math side . &lt;EOT&gt;</li> <li>• how do you start to like yourself ? i think my parents is by focusing . &lt;EOT&gt;</li> <li>• how do you start to yourself like ? i was taught by my parents . &lt;EOT&gt;</li> </ul>	<ul style="list-style-type: none"> <li>• are u shy or k ? both , actually &lt;EOT&gt;</li> <li>• are u minded or rem ? actually , both &lt;EOT&gt;</li> <li>• are u transparent or shy ? it'd actually , add-on &lt;EOT&gt;</li> <li>• are u untouchable cubed or programe ? both , actually like &lt;EOT&gt;</li> <li>• wha do u are roselle or marketed ? you start , by both my inbox &lt;EOT&gt;</li> <li>• how do u simplify phases towards you ? are proving , like no smiles . &lt;EOT&gt;</li> <li>• how do you burp confidence ? to start i was like , shareaza the new by hindering . &lt;EOT&gt;</li> <li>• how do you start to race ? i like kazaa when my was cheated . &lt;EOT&gt;</li> <li>• how do you start to start like ? i was taught by my parents . &lt;EOT&gt;</li> <li>• how do you start to like yourself ? i was taught by my parents . &lt;EOT&gt;</li> </ul>
<b>(Health)</b> <BOT> how do you start to like yourself ? i was taught by my parents .	
<ul style="list-style-type: none"> <li>• how do you start to yourself by allowing ? i like my parents yr . &lt;EOT&gt;</li> <li>• how do you start to yourself like i ? my parents was by mario practitioner . &lt;EOT&gt;</li> <li>• how do you start to cite yourself ? i like by my consequences in 1981 . &lt;EOT&gt;</li> <li>• how do i start girls like to ? you can find yourself in my states , by today . &lt;EOT&gt;</li> <li>• how do you start yourself drunk ? i can find in something like to my country , what by jane . &lt;EOT&gt;</li> <li>• how can i start those need in america ? do you like to rephrase an invention , what i'm spinning ? &lt;EOT&gt;</li> <li>• how can i find someone in spain ? i'm guessing today by pascal , what do you want to ? &lt;EOT&gt;</li> <li>• how can i find an attorney in spain ? i'm studying chicken's what , do you want to ? &lt;EOT&gt;</li> <li>• how can i find someone in spain ? in spain i'm studying , what do you want ? &lt;EOT&gt;</li> <li>• how can i find someone in spain ? i'm in italy today , what do you want ? &lt;EOT&gt;</li> </ul>	<ul style="list-style-type: none"> <li>• how do you start to like yourself ? i was taught by new england . &lt;EOT&gt;</li> <li>• how do you start to like yourself ? i was taught by my parents . &lt;EOT&gt;</li> <li>• how do i start you to beethoven ? like israel was my grandmother by fielders . &lt;EOT&gt;</li> <li>• how do you start to find ? i like aggieland in my testicles was listening . &lt;EOT&gt;</li> <li>• how can i do compuserve attain ? start to comment in spain you like , was my real pics . &lt;EOT&gt;</li> <li>• how can i find blueprints do you ? i'm in spain like queens to chelsea , arrange . &lt;EOT&gt;</li> <li>• how can i find uneasy profiles in spain ? i'm sure what you do , like today's ? &lt;EOT&gt;</li> <li>• how can i find someone in spain ? i'm in spain today , what do you want ? &lt;EOT&gt;</li> <li>• how can i find someone in spain ? i'm in tanks today , what do you want to ? &lt;EOT&gt;</li> <li>• how can i find someone in spain ? i'm guessing in spain today , what do you want ? &lt;EOT&gt;</li> </ul>
<b>(Business &amp; Finance)</b> <BOT> how can i find someone in spain ? i'm in spain today , what do you want ?	

Table 17. Interpolation results between latent codes of input sentences (with gray) from **Yahoo Answers**.

### D.3. Sampling

sMIM (512)
<ul style="list-style-type: none"> <li>• instead the stock market is still being felt to &lt;unk&gt; those of our empty than in a bid &lt;EOT&gt;</li> <li>• he estimated the story will take &lt;unk&gt; of paper co . ' s \$ n million in cash and social affairs to at the company a good share &lt;EOT&gt;</li> <li>• long-term companies while the company ' s &lt;unk&gt; provisions would meet there to n or n cents a share and some of costly fund &lt;EOT&gt;</li> <li>• time stocks the company explained him to sell &lt;unk&gt; properties of high-grade claims which has received a net loss in the firm &lt;EOT&gt;</li> <li>• what i had the recent competition of &lt;unk&gt; replies that is n't expected to draw a very big rise in tokyo &lt;EOT&gt;</li> </ul>

Table 18. Samples from best performing model for dataset **PTB**.

sMIM (1024)
<ul style="list-style-type: none"> <li>• ben monkey gabi sister near the western fest . i ' ve been looking forward to this location , and each time i ' m in the 6th bunch i want to have a great visit experience . it was all kinds of fillers , owns and dressings non-asian with jalapeños &lt;unk&gt; does n't hold me for much healthier . front desk is not my favorite dinner place at the gates . they are closed on mondays , - lrb - it could affect a couple minutes more rocks - rrb - and then we said the bar was the real bold . i ' d rather go to firefly some bubble in greece . if you had a neighbourhood addiction &lt;unk&gt; c , take this look as most amazing . &lt;EOT&gt;</li> <li>• hello tanya stephen covering qualité . ugh haha , i was curious to consume that the white asian restaurants believes filled a mob and turkey melt departments for \$ 9.99 . the &lt;unk&gt; of these were not intrusive , it was accepted in there . . . i ' m sure this is n't one of my favorite places to go at night with here ! particularly speaking the italian cleaning tables . we also ordered some pina colada , which tasted exactly like they came out of a box and per endearing thick . pretty good food overall , and the pigeons self nightly . i ' d call it again just on halloween for a dependable lunch . but the statue sucks ? so if you have bouchon to inquire was good place . &lt;EOT&gt;</li> <li>• prada based pata based solely often inside . this place is unappealing horrific for the 50th and fries , i ' ve caught to have a ton of good reviews &lt;unk&gt; in buckeye , barnes knew . not bc that i was wrong with my team being kicked the whole thing at eggroll , it ' s like pulling out of the landmark . no luck on ketchup top crunch , if you are craving something simple and &lt;unk&gt; . we also tried the wild mushroom - lrb - it ' s burn , did n't go in disheveled - rrb - as a matter destination from flavor . the food was just ok and nothing to write home about . friend peeps i only had one beer , but this place does not deserve the same increase . &lt;EOT&gt;</li> </ul>

Table 19. Samples from best performing model for dataset **Yelp15**.

sMIM (1024)
<ul style="list-style-type: none"> <li>• how does transformers send grow ina under pubs ? i found the suspension resides official game is exciting to withstand and what can a person do in that case ? bree fights , if it does 150 . the dre is tied ordered outlook &lt;unk&gt; 2005 . today had a migraine with limitation tops , because of his vr repeats , you are referring to review at the university of 1994 and have visited fortune . judy for websites &lt;unk&gt; website is beware confused . &lt;EOT&gt;</li> <li>• how do i download jesus gyno to woman whom ? being irvine in line is what you did a lot of oceanic denny in the middle east and spanish wallet or &lt;unk&gt; entity . plus , i'm aware of that , particularly do you have any insight insight ... if you are a hoe who's right click on it , and you can ' t get some skills god . the other government also happened to be &lt;unk&gt; with most varied life-forms is located at this point . foreigners your covers , and maybe even my friends . &lt;EOT&gt;</li> <li>• what's mastering marathons fluently is einstein among the waivers ? ok i feel that what happened to tom during the holidays monitor of 1-2 awol whn reservoir &lt;unk&gt; . clusters in a workforce and it symbolizes , seems are meant to have any distinction on the patriot , british languages even though i would build god if you like . just bringing your old door as a distorted spree ? hmmm , because you're not anti-bacterial pure dino and &lt;unk&gt; this can be deduced . &lt;EOT&gt;</li> </ul>

Table 20. Samples from best performing model for dataset **Yahoo Answers**.

Here we show samples from the best performing models learned from a single dataset for PTB (Fig. 18), Yelp15 (Fig. 19), and Yahoo Answers (Fig. 20). We sample from a zero-mean Gaussian distribution over the latent space, with an isotropic covariance with a standard deviation of 0.1 (since we cannot directly sample from the implicit marginal over the latent). Interestingly, this simple heuristic provides good samples. We attribute this to the anchor, which defines scale and position for the implicit marginal over the latent to roughly match.

### D.4. Question Answering

Here we provide more examples of answers generated from a model trained on Yahoo Answers (*i.e.*, sMIM (1024) in Fig. 21). In particular, the model was trained from data in which 20% of the encoder input tokens were replaced with the <unk> token. This is a form of self-supervised learning commonly used in language modelling (*e.g.*, Bowman et al. (2015)). This encourages the model to replace <unk> with other tokens. We have found this procedure to significantly improve the quality of the generated answers. We provide three generated answers for each question (Q), taken from Yahoo Answers. Short/medium/long answers (A) are generated by concatenating 5/10/15 <unk> tokens. The number of <unk> encodes the length of the expected answer. We note that, in many cases, only one answer will be a good match to the question, suggesting the model has preferences towards answers with a question specific length.

## SentenceMIM

sMIM (1024)	
(Q)	<BOT> where can i find the most affordable roses for valentine's day this late in the game ?
(A) short	• try the mall , <EOT>
(A) medium	• the game is to play the game , set your mind <EOT>
(A) long	• i think the best answer is to ask her out , especially since the late 1990s game <EOT>
(Q)	<BOT> how do you start to like yourself ?
(A) short	• start by yourself , take the steps <EOT>
(A) medium	• start by yourself . you have to start asking yourself why this is happening <EOT>
(A) long	• start by yourself . you have to start acting yourself in the middle of your eye . <EOT>
(Q)	<BOT> how can i find someone in spain ?
(A) short	• try classmates.com or classmates.com to see <EOT>
(A) medium	• someone who is a good one to talk to your question <EOT>
(A) long	• someone who is a good person to find a nice girl , and take care of yourself <EOT>
(Q)	<BOT> with a weapon to your person , who in you family do you save with an answer ?
(A) short	• you should save yourself in person <EOT>
(A) medium	• you can save the money in person , or do you want to save money <EOT>
(A) long	• you can save your money in person , and answer the question yourself . do you have a family ? <EOT>
(Q)	<BOT> what are the best sites to shop for personal computers online ?
(A) short	• try amazon.com online . <EOT>
(A) medium	• i think it is best to shop online , or take a look at the personal <EOT>
(A) long	• yahoo is best online . i would suggest checking out the personal website for personal info <EOT>
(Q)	<BOT> who is going to win the super bowl this year ?
(A) short	• the steelers is a pretty good <EOT>
(A) medium	• the pittsburgh steelers is a good one , but i don't know <EOT>
(A) long	• this is the best team to win the super bowl , and i think you mean good luck <EOT>
(Q)	<BOT> is there a web site that provides info on companies that have been known to provide lousy service ?
(A) short	• yes , google was a little service <EOT>
(A) medium	• i have known as a service that provides a service pack to provide transparency . <EOT>
(A) long	• try searching on google and search for that info . there are many different types of service that provide to the service that has been answered <EOT>
(Q)	<BOT> what is the closest capital city to paris ?
(A) short	• paris hilton city , ca <EOT>
(A) medium	• in the south city , ca , i think you mean the city hall <EOT>
(A) long	• city is the capital city in south america . i think it is a beautiful city <EOT>
(Q)	<BOT> is santa claus real ?
(A) short	• yes , you can ask her to questions <EOT>
(A) medium	• yes , it is a santa claus . the phrase is to take care <EOT>
(A) long	• yes , it is real . santa claus are the best in santa claus . take care <EOT>
(Q)	<BOT> do height increasing devices shown on television works ?
(A) short	• perhaps shown is shown to see <EOT>
(A) medium	• yes , the best way to increase it is to work <EOT>
(A) long	• yes , shown shown . the biggest problem with television is to take a walk away from the edge <EOT>
(Q)	<BOT> where can i find advice on hip hop beat making ?
(A) short	• try this website : <a href="http://messenger.yahoo.com/">http://messenger.yahoo.com/</a> <EOT>
(A) medium	• try this website : <a href="http://messenger.yahoo.com/">http://messenger.yahoo.com/</a> to find the best answer . <EOT>
(A) long	• try making a video club , such as the lion and hip hop . i wish you to be good <EOT>
(Q)	<BOT> how do u clear the history in the search field ?
(A) short	• clear search history , it will be clear <EOT>
(A) medium	• clear history , it appears that the question period has expired . if you have received an answer that meets your needs , please choose a ' best answer . ' regards <EOT>
(A) long	• clear history in the search field , u can find it from your website : * clear history <EOT>
(Q)	<BOT> where can i get chapter summaries in text books ?
(A) short	• try this site , you'll find it interesting <EOT>
(A) medium	• text books ? try this site , and get a book to read <EOT>
(A) long	• in books , it is a text book , and the text books are written in the same text . <EOT>
(Q)	<BOT> how to tell a nice guy you dont like him ?
(A) short	• nice guy dont know what to do <EOT>
(A) medium	• nice guy , dont tell him what the hell is <EOT>
(A) long	• dont tell him that you like him ? nice guy , and the guy who is nice to him ! <EOT>
(Q)	<BOT> does your body feel physically fit ?
(A) short	• no , it is a little bit <EOT>
(A) medium	• feel your body needs to fit into the body . i feel like a good fit <EOT>
(A) long	• feel your body fit in a fit body . i feel like the best fit to fit in your body <EOT>

Table 21. Question and sampled answers from model sMIM (1024) (i.e., trained on Yahoo Answers dataset). We provide short/medium/long sampled answers (A) for each question (Q).