# Regularizing Transformers With Deep Probabilistic Layers

Aurora Cobo Aguilera[a,*], Pablo Martínez Olmos[a], Antonio Artés-Rodríguez[a], Fernando Pérez-Cruz[b]

[a]*Department of Signal Theory and Communications*
*Universidad Carlos III de Madrid*
*Avda. de la Universidad 30, 28911, Leganés, Madrid, Spain*
[b]*Swiss Data Science Institute (ETHZ/EPFL)*
*Universitatstrasse 25, 8006, Zurich, Switzerland*

## Abstract

Language models (LM) have grown with non-stop in the last decade, from sequence-to-sequence architectures to the state-of-the-art and utter attention-based Transformers. In this work, we demonstrate how the inclusion of deep generative models within BERT can bring more versatile models, able to impute missing/noisy words with richer text or even improve BLEU score. More precisely, we use a Gaussian Mixture Variational Autoencoder (GMVAE) as a regularizer layer and prove its effectiveness not only in Transformers but also in the most relevant encoder-decoder based LM, seq2seq with and without attention.

*Keywords:* Natural Language Processing, Regularization, Deep Learning, Transformers, Variational Auto-Encoder, missing data

## 1. Introduction

Deep Generative Models (DGMs) have become a cornerstone in modern machine learning due to their ability to learn abstract features from high-dimensional spaces to generate new data (Goodfellow et al., 2014, Kingma and Welling,

---

*Corresponding author

*Email addresses:* `acobo@tsc.uc3m.es` (Aurora Cobo Aguilera), `pamartin@ing.uc3m.es` (Pablo Martínez Olmos), `aartes@ing.uc3m.es` (Antonio Artés-Rodríguez), `fernando.perezcruz@sdsc.ethz.ch` (Fernando Pérez-Cruz)

2014). In the field of Natural Language Understanding (NLU), state-of-the-art is dominated by attention-based probabilistic models, a class of explicit DGMs that can be trained with Maximum Likelihood Estimation (MLE) approaches (Caccia et al., 2020).

Regarding other well known DGMs such as Generative Adversarial Networks or GANs (Goodfellow et al., 2014), so far for NLU they have not shown the same outstanding results that they achieve for image processing (Zhang et al., 2019, Radford et al., 2015), mostly due to the discrete nature of the data, which leads to non-differentiable issues, mode collapse and optimization instability (Lu et al., 2018, Caccia et al., 2020). To tackle these and other issues, recent contributions propose the use of Reinforcement Learning techniques to optimize the GAN loss function (Yu et al., 2017, Fedus et al., 2018, Guo et al., 2018, de Masson d'Autume et al., 2019), continuous approximations to discrete sampling (Jang et al., 2017, Zhang et al., 2017), or learning a low-dimensional representation through autoencoders (Zhao et al., 2018, Subramanian et al., 2018, Donahue and Rumshisky, 2018, Yu et al., 2018, Haidar et al., 2019, Haidar and Rezagholizadeh, 2019, Rashid et al., 2019). Besides, explicit DGMs such as variational autoencoders (VAEs) have also been proposed in several NLU approaches again with limited success (Pagnoni et al., 2018, Shen et al., 2018, Gupta et al., 2018, Yang et al., 2017, Shi et al., 2019). Some of the pioneers in this field were Bowman et al. (2016), who proposes a RNN-based VAE for text generation. Even in an extent, Hu et al. (2017) combine a VAE with a discriminator to build a hybrid model that solves the text generation problem. In all these works, both GANs and VAEs are at the core of the NLU model, and hence are fully responsible to capture the semantic structure and generate text. For this particular task, they are still not competitive with attention-based probabilistic models (Caccia et al., 2020).

In this work, we propose to exploit DGMs for NLU in a completely novel and different way. Instead of training a DGM to solve a NLU task, we rely on a hybrid model in which a transformer-based architecture like BERT (Devlin et al., 2018) is combined with a VAE, which is placed inside its structure as a stochastic

layer that helps to learn a richer hidden space, enforcing a regularization effect to some extent. In particular, we use a structured VAE that implements a mixture of Gaussians in the latent space (GMVAE) (Dilokthanakul et al., 2016), since it is able to capture more complex data in an easier way than the traditional vanilla VAE. In a similar way, Sriram et al. (2018) and Gulcehre et al. (2015) built fusion models taking advantage of a pre-training process as we explain later. Nevertheless, they only focused on a basic seq2seq architecture.

Regularization in deep learning has risen up from the beginning of Neural Networks with the extensively use of tools such as dropout (Srivastava et al., 2014), early stopping, data augmentation or weight decay (Krogh and Hertz, 1992), which helps models to generalize. However, regularization in NLUs is a much-less explored field and none of these tools experience the same versatility as our proposal in this paper, in which the GMVAE performs a controlled and structured noise injection within the NLU deep network. When combined with BERT, we name our model as NoRBERT (Noisy Regularized BERT) and we conclude that the effect of the stochastic layer is very different depending on the transformer layer where it is placed. If the layer is placed at the end of the structure, it drives more versatile topics when imputing missing words. On the contrary, when placed at the bottom, it improves BLEU score, what coincides with the goal of traditional regularization mechanisms.

We illustrate our approach in word imputation problems (masking the source text corpora) using a BERT transformer network, demonstrating gains in machine translation setups (better BLUE scores) and the versatility of the method to impute missing words by a large set of examples. Furthermore, we also explore the GMVAE regularization effect in traditional seq2seq models with and without attention mechanisms and explain the regularization functionality in a simple well-known problem as it is classification of Fashion MNIST images. The code to generate our results is available in an open repository[1].

This paper is organized as follows. Firstly, in Section 2 we describe some

---

[1]https://github.com/AuroraCoboAguilera/NoRClassifier

related work which is key to understand the paper: the VAE, and more precisely the GMVAE, as the main structure of the regularizer and transformer networks with BERT as our model to be studied. Secondly, in Section 3 we explore a basic example of applying our idea in a well-known scenario as it is Fashion MNIST. This is a useful prove of the stochastic layer effect and its effectiveness in other problems. Thirdly, in Section 4 we describe in detail our model, NoRBERT, and two variants of it, Top and Deep NoRBERT, depending on the transformer layers where we apply the regularization. Then, we present the results of these two options in Section 5. Moreover, we include an extension (Section 6) where we study other relevant encoder-decoder based LMs as it is seq2seq with and without attention (Bahdanau et al., 2015). Finally, in Section 7 we conclude our work and mention some future lines of research.

## 2. Related work

### 2.1. Variational Autoencoders with Gaussian mixture priors

A VAE (Kingma and Welling, 2014) is a class of density estimator that consists on two networks, an encoder and a decoder or generator, that builds a regular latent space with the help of probability distributions. The properties of the organized latent space allow not only the reconstruction of the input data but also the generation of new instances from a sampling procedure. In a standard vanilla VAE, see Figure 1a, the low-dimensional latent space follows a Gaussian prior distribution likelihood parameters, e.g. mean and covariance matrix of $p(x|z)$ are parameterized with the decoder network with input $x$. Variational inference of the model parameters is achieved by maximizing a lower bound on $\log p(x)$, which in turn depends on a flexible NN parameterized distribution $q(z|x)$ that approximates the true posterior $p(z|x)$:

$$\mathcal{L}_{ELBO}\left(\theta, \phi, x\right) = \mathbb{E}_{z \sim q_\phi(z|x)}\left[\log p_\theta\left(x|z\right)\right] - \mathcal{KL}\left[q_\phi\left(z|x\right) \| p\left(z\right)\right], \qquad (1)$$

where $\mathcal{KL}(q|p)$ is the KL divergence between distributions $q$ and $p$ and acts as a regularization in the evidence lower bound (ELBO) objective. The graphical model of $q(z|x)$ is indicated in Figure 1a with dotted lines.
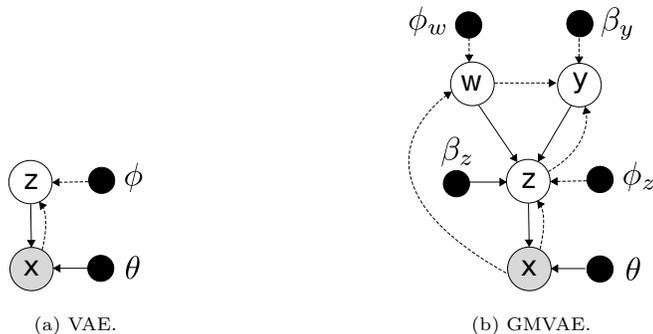
(a) VAE.  (b) GMVAE.

Figure 1: The directed graphical models into consideration. Solid lines denote the generative model and dashed lines the variational approximation. The shaded variables are considered the observed inputs, the dark units are the networks parameters to be optimized and the units that are left are the latent variables.

The flexibility of VAEs have encouraged the study of different priors and architectures to obtain models capable of inferring more complex structured data. That is the case of using a Mixture of Gaussians (MoG) as the prior distribution $p(z)$ for the latent space because it helps to capture the multimodal nature of some data (Jiang et al., 2017, Dilokthanakul et al., 2016). We refer to this model as GMVAE, and its graphical model is shown in Figure 1b. The GMVAE generative model proposed by Dilokthanakul et al. (2016) is characterized by the following distributions:

$$p(z) = \int p(z|w,y) \cdot p(w) \cdot p(y) dw dy \tag{2a}$$

$$p(w) = \mathcal{N}(0, \mathcal{I}) \tag{2b}$$

$$p(y) = \text{Mult}(\pi), \quad \pi_i = \frac{1}{K} \tag{2c}$$

$$p_{\beta_z}(z|w,y) = \prod_{k=1}^{K} \mathcal{N}(\mu_{\beta_{y_k}}(w), \Sigma_{\beta_{y_k}}(w))^{y_k == 1} \tag{2d}$$

$$p_\theta(x|z) = \mathcal{N}(\mu_\theta(z), \sigma\mathcal{I}), \tag{2e}$$

where $\mu_{\beta_{y_k}}$, $\Sigma_{\beta_{y_k}}$ and $\mu_\theta$ are neural networks. $\mu_{\beta_{y_k}}$ and $\Sigma_{\beta_{y_k}}$ indicate a different NN per component in the mixture of Gaussians and $K$ is the total number of components. The posterior distribution of $z$, $w$ and $y$ given $x$ is chosen according

to the following factorization

$$q_{\phi_z}(z|x) = \mathcal{N}(\mu_{\phi_z}(x), \Sigma_{\phi_z}(x)) \tag{3a}$$

$$q_{\phi_w}(w|x) = \mathcal{N}(\mu_{\phi_w}(x), \Sigma_{\phi_w}(x)) \tag{3b}$$

$$q_{\beta_y}(y_j == 1|w, z) = \frac{p(y_j == 1) \cdot p_{\beta_z}(z|y_j = 1, w)}{\sum_{k=1}^{K} p(y_k == 1) \cdot p_{\beta_z}(z|y_k = 1, w)}, \tag{3c}$$

where again $\mu_{\phi_z}$, $\Sigma_{\phi_z}$, $\mu_{\phi_w}$, and $\Sigma_{\phi_w}$ are dense neural networks, resulting in the following evidence lower bound (ELBO):

$$\mathcal{L}_{ELBO}(\theta, \phi, x) =$$
$$\mathbb{E}_{z \sim q_{\phi_z}}[\log p_\theta(x|z)] - \mathbb{E}_{w \sim q_{\phi_w},\, y \sim p_{\beta_y}}[\mathcal{KL}[q_{\phi_z}(z|x) \| p_{\beta_z}(z|w, y)]] - \tag{4}$$
$$\mathbb{E}_{z \sim q_{\phi_z},\, w \sim q_{\phi_w}}[\mathcal{KL}[p_{\beta_y}(y|w, z) \| p(y)]] - \mathcal{KL}[q_{\phi_w}(w|x) \| p(w)]$$

### 2.2. Transformer networks: BERT

Over the last couple of years, Transformers (Vaswani et al., 2017) have become a revolution in the field of NLU (Dai et al., 2019, Keskar et al., 2019, Ma et al., 2019, Gu et al., 2019, Yang et al., 2020b) due to their ability to capture longer-range linguistic structure. Unlike previous works (Sutskever et al., 2014, Bahdanau et al., 2015, Luong et al., 2015), they rely entirely on self-attention to compute the latent representations of the sentences.

Transformer-based models are usually applied in a transfer learning perspective (Devlin et al., 2018, Radford et al., 2018, Song et al., 2019, Radford et al., 2019, Yang et al., 2019, Lample and Conneau, 2019, Dong et al., 2019, Sun et al., 2020, Xiao, 2020, Yang et al., 2020a) that allows users to train smaller datasets in a specific task quicker and more accurate than doing it from scratch. Firstly, you need a pre-trained model that has learned contextualized text representations in a general unsupervised scenario with a large text corpus. Afterwards, you can fine-tune the model using a small database with the addition of few parameters or layers in a downstream task. This is the case of BERT (Devlin et al., 2018), which stands out above all, providing a pre-trained Transformer text encoder as a general LM for any downstream task. Since its appearance, several

BERT-based models have emerged (Liu et al., 2019, Sanh et al., 2019) and to-day they dominate the leaderboard[2] in GLUE benchmarks (Wang et al., 2019a).
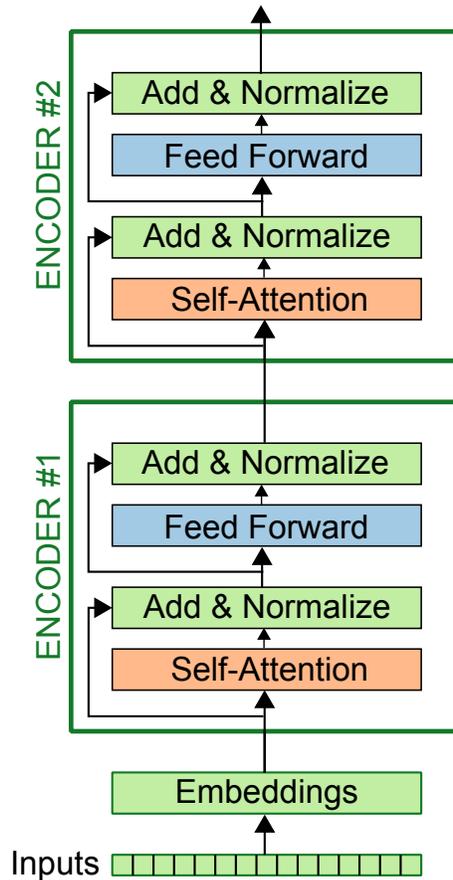


Figure 2: Diagram of BERT structure.

Figure 2 shows a diagram with the structure of BERT for the first two layers. Basically, it is composed of a first step with the computation of the input sentences embeddings, then a pile of transformer encoder layers, and then, if it is necessary we can apply any task-specific layer on top. Each of these layers consists on two blocks, a multi-head self-attention mechanism and a feed forward network, both with a normalization following them. Regarding the im-

---

plicit regularization mechanisms within BERT, dropout and weight decay are applied through all the structure: in the fully connected layers in the embeddings, encoder, pooler and in the attention probabilities with rates of 0.1 a 0.01 respectively.

BERT makes use of WordPiece embbedings (Wu et al., 2016) with a vocabulary size of 30000 tokens. The base models are pre-trained in the datasets of Book Corpus (Zhu et al., 2015) with 800M words and English Wikipedia with 2500M words.

Although we focus our work in NoRBERT, we will extend the results to traditional seq2seq models (Section 6) in order to explain how our mechanism works and show its ability to be integrated in other architectures.

## 3. GMVAE as a regularizer in deep neural networks

In this work, we put forward GMVAEs as a robust stochastic layer to enforce regularization in a deep NN, with particular focus on Transformers and NLU. Before describing the methodology in a complex transformer based network, we want to illustrate our approach in a simpler setup, in which we regularize a deep six-layer MLP over the Fashion MNIST (FMNIST) database[3] ([dataset] Xiao et al., 2017).
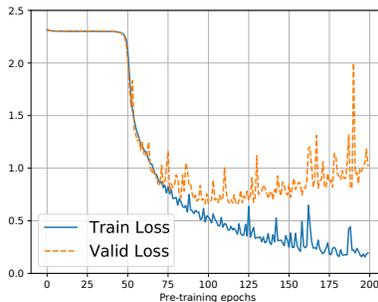


Figure 3: Pre-training a six layer MLP for FMNIST.

---

8

In Figure 3 we show the train/validation cross entropy loss of the NN in a completely unregularized training (no dropout or weight decay whatsoever). Validation error begins to raise up from epoch 120. Now we perform the following experiment. We get the NN parameters at epoch 119, at which overfitting was not yet noticeable, and we introduce two types of regularization layers between the first two MLP layers:

1. A standard dropout layer with erase probability $p$.
2. A GMVAE layer trained using the 700-dimensional internal representation of the first MLP layer. For every output from the first MLP layer, the GMVAE layer first computes a latent low-dimensional representation sampling from the GMVAE posterior distribution in (3a)-(3c) to then provide at the output a reconstruction sampled from generative model in (2a)-(2e).

The details about the GMVAE layer parameters used for this experiment can be found in Appendix C.1. Note that the GMVAE layer, as dropout, is introducing a certain level of distortion over the input vector but, unlike dropout, such distortion is not independent to the input vector, as for some atypical vectors the reconstruction noise will be larger. This allows the network to explore diverse regions at the input of the following layer. In Figure 4 we show the train/validation cross entropy loss when layer 1 is fixed (so the GMVAE input distribution is not changing) and we keep training MLP layers 2-6. In Figure 5 we show the performance when dropout with $p = 0.1$ (a) and $p = 0.5$ (b) is used instead of the GMVAE layer. On the one hand, observe the inability of dropout to compensate the overfitting of the network. On the other hand, due to the controlled noise injection, the GMVAE avoids overfitting even after an excess of additional epochs. With these figures we can state that the training loss decays much more slowly in our model with a score of 0.3 after 700 epochs, while in the dropout case it drops off almost to zero after 150 epochs.

With this example, we simply want to put forward the use of a DGM (a GMVAE in our case) as potential regularizer with additional flexibility, compared to simpler solutions such as dropout. A detailed cross-validation analysis of what

9

Figure 4: Fine-tuning a six layer MLP for FMNIST with the GMVAE regularization layer placed after the first MLP layer.

kind of regularization method optimizes the classification performance in this particular classification setting is not relevant at this point. In the following, we show how the use of GMVAE layers if able to enhance the performance of complex pre-trained networks such as BERT, which of course has already been trained with its own regularization methods (including dropout).



(a) With dropout probability of 0.1.

(b) With dropout probability of 0.5.

Figure 5: Fine-tuning a six layer MLP for FMNIST with dropout in the first layer.

## 4. Improving BERT with GMVAE layers: NoRBERT

### 4.1. Overview

The main idea of our work is the integration of the GMVAE in BERT through NoRBERT. In this hybrid model, the GMVAE layer alters the BERT hidden embeddings in one particular layer through a project-and-reconstruct operation, adding a structured noise to them and hence enforcing a regularization mechanism. In other words, we try to break the determinism in exchange of more robust solutions. Unlike in other regularization techniques such as dropout, the reconstruction error plus the observation noise (GMVAE noise for short) of the GMVAE will not be uniform across embeddings, since atypical embeddings will suffer from larger GMVAE noise variance. As a result, the network training will rely less on such noisy embeddings, which we show is beneficial for the overall performance.

We want to stress the fact that we use BERT as an exemplary case of how a certain neural language model can be enhanced by the inclusion of GMVAE layers within. Furthermore, in Section 6 we show how to incorporate the same idea in seq2seq language models with attention. Moving back to BERT, NoRBERT builds upon a pre-trained BERT model, allowing the integration of the GMVAE in an intermediate step. We follow these four main steps:

1. Pre-train BERT with a masked text corpora.
2. Train a GMVAE over the space of hidden embeddings coming from input sentences using one particular BERT layer.
3. Include the GMVAE layer inside the structure. The GMVAE will be responsible for adding noise in the propagation of the information, as in the GMVAE layer every input vector is projected into a low-dimensional space and reconstructed back by sampling from the generative model.
4. Retrain the model by fine-tuning all layers above the GMVAE one. The layers below the GMVAE one are not altered so we do not modify the embedding space in which the GMVAE was trained on.

11

Regarding the base BERT model, for the implementation we use the base model from Devlin et al. (2018). In the training we use the masked language modeling (MLM) strategy as Liu et al. (2019), since it is the straightforward strategy to train transformers in word imputation (Song et al., 2019).

### 4.2. Top and Deep NoRBERT

In the study of NoRBERT we explore placing the regularizer in different layers from BERT. Firstly, we show the effect of the GMVAE on top of the transformer encoder, just before the classification layer that computes the vocabulary logits. We refer to this case as *Top NoRBERT*. Secondly, we explore the consequences when the biggest part of BERT is retrained after placing the GMVAE in one of the first and middle layers. This is referred to *Deep NoRBERT*.
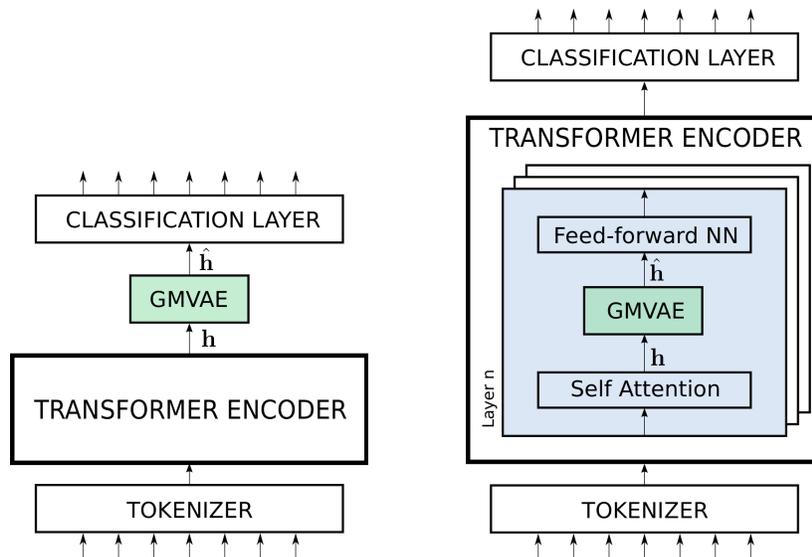


Figure 6: Top NoRBERT (a) and Deep NoRBERT (b).

Top NoRBERT consists on a new version of BERT, with a stochastic layer on the top of the transformer encoder as represented in Figure 6(a). Therefore, the only difference from the original model is that we use a GMVAE to reconstruct the last hidden states before the final token decision. We previously train

the GMVAE with the hidden states computed by base BERT for the training sentences. Afterwards, we fine-tune the classification layer of BERT with the stochastic reconstruction integrated.

In Deep NoRBERT we include the GMVAE stochastic layer inside an intermediate transformer layer, after the self-attention and before the feed-forward blocks as shown in Figure 6(b). We fine-tune the parameters in the structure above the regularizer, that is, the feed-forward block in the same encoder layer and the transformer layers that are on top of it. In our experimental results, we demonstrate gains w.r.t. the base BERT model by including only one GMVAE layer. We explore the method in several layers with results that are qualitatively very different compared to Top NoRBERT.

## 5. Results

To implement NoRBERT, we make use of the pre-trained base model from BERT described by Devlin et al. (2018). This version of BERT is composed of 12 layers, a hidden size of 768 and 12 heads and we make use of the parameters eased by the *Hugging face* library[4] using a MLM objective. We keep the original configuration following the paper (Devlin et al., 2018) except for the hyperparameters mentioned in the following sections. On each experiment, we train a GMVAE using the hidden vectors at some point of BERT structure obtained from training samples with this base model. Once the GMVAE has converged, we build a new architecture based on BERT with the integration of a stochastic layer in the corresponding place of the hidden vectors. This new layer consists on the reconstruction of the hidden vectors through the generative network of the GMVAE. Finally, we fine-tune this new architecture, freezing all parameters below the stochastic layer in the computational graph.

In the experiments we employ the dataset *snli*[5] ([dataset] Bowman et al., 2015) with different strategies in the masking process of tokens. It has a vo-

---

[4]https://huggingface.co/
[5]https://nlp.stanford.edu/projects/snli/

cabulary size of 36711 different words. We use the entire preprocessed training data which contains 714667 sentences and a test set of 13350. The GMVAE is trained for all the tokens of a random set of 50000 training sentences. When saving the hidden states to train the GMVAE a posteriori, we treat each token as an independent input to the GMVAE, ignoring tokens that correspond to padding (they exist due to BERT format of the tokenizer).

To speed up the loading of data, we utilize the extension *hdf5* for saving the files. Moreover, this way we avoid memory issues when loading the datasets in the programs since these files with the hidden states have significant sizes.

### 5.1. Deep NoRBERT

First, we present the results of Deep NoRBERT, in which the GMVAE stochastic layer is placed in an intermediate BERT encoder layer, see Section 4.2. Next we present the results obtained in terms of accuracy and BLEU score for different locations of the GMVAE layer inside the BERT structure.

The GMVAE layer is trained for 500 epochs with a learning rate of $5 \cdot 10^{-5}$. The GMVAE latent dimension $z$ is set to 150, the $w$ dimension to 50, and we consider a mixture of 20 Gaussians, dropout 0.3 and networks with a depth of 6 layers. Then, deep NoRBERT is trained for 8 epochs freezing the parameters below the stochastic layer. The baseline BERT is also fine-tuned in the same dataset for 8 epochs so we can make a fair comparison in their performance in missing data imputation. We evaluate the percentage of tokens that are exactly the same as the source sentence in a 1-by-1 comparison. We test two different scenarios, with masked tokens and with disrupted tokens, that is, instead of using the [MASK] token which indicates 'unknown', we place random choices from the vocabulary that damage the source sentence. We replicate the random words substituted on each experiment maintaining the same seed in the training. Regarding the masks, 40% of the sentences chosen at random have at least one [MASK] token, which always replaces a meaningful word (we avoid masks over stopping words).

Table 1 shows the imputation accuracy for different configurations, in which

14

| Model | Masked | Disrupted |
|---|---|---|
| BERT | 97.13% | 96.98% |
| 1-Deep NoRBERT | **97.32%** | **97.11%** |
| 2-Deep NoRBERT | **97.20%** | **97.07%** |
| 3-Deep NoRBERT | **97.18%** | **97.1%** |
| 9-Deep NoRBERT | 96.87% | 96.25% |
| 11-Deep NoRBERT | 96.05% | 95.34% |
| 12-Deep NoRBERT | 95.89% | 93.89% |

Table 1: Accuracy of different models comparing the unmasked source sentence with the reconstruction. We evaluate a version that keeps the [MASK] tokens and other that substitutes them by random tokens from the vocabulary. In $l$-Deep Norbert, $l$ refers to the transformer BERT layer in which the GMVAE is placed.

$l$-Deep NoRBERT means that we placed the GMVAE layer in the $l$-th transformer layer. For a better visualization, we highlight in bold every case that outperforms the baseline. Observe that the largest gains are obtained when the GMVAE layer is placed in the bottom of the network, outperforming BERT after fine-tuning. We remark that BERT is a state-of-the-art model for MLU that is pre-trained over a massive dataset and hence any improvement is not negligible, particularly when is achieved by placing a single regularization layer within. Despite some studies about BERT state that the last layers encode task-specific features (Kovaleva et al., 2019), our results demonstrate that fine-tuning and regularization of deep layers may improve the overall performance.

Table 2 presents the BLEU score obtained by Deep NoRBERT with different layer configurations. We explore different policies of generating missing tokens. 'Low' refers to the same mechanism as in Table 1 experiments. In the policies called 'Medium' and 'High' we do not exclude any token by its grammatical meaning and mask every word independently with probabilities of 0.4 and 0.6 respectively. Table 3 results, called Masked BLEU, differ from the previous ones in the n-grams taken for the metric computation. That is, we only consider n-grams that include a masked token. From both tables we draw similar

| Model/Missing rate | Low | Medium | High |
| --- | --- | --- | --- |
| BERT | 86.07 | 49.43 | 25.14 |
| 1-Deep NoRBERT | **86.90** | **49.91** | **25.53** |
| 2-Deep NoRBERT | **86.65** | **49.75** | **25.26** |
| 3-Deep NoRBERT | **86.53** | 49.33 | **25.45** |
| 9-Deep NoRBERT | 85.52 | 46.04 | 21.47 |
| 11-Deep NoRBERT | 83.89 | 43.34 | 19.28 |
| 12-Deep NoRBERT | 80.77 | 40.83 | 17.16 |

Table 2: BLEU score of different models comparing different missing rates.

conclusions: the best performance is obtained when the GMVAE layer is placed at the bottom of the network, right after the first transformer layer.

| Model/Missing rate | Low | Medium | High |
| --- | --- | --- | --- |
| BERT | 3.73 | 21.3 | 15.34 |
| 1-Deep NoRBERT | **3.88** | **22.7** | **16.44** |
| 2-Deep NoRBERT | **3.88** | **22.50** | **16.22** |
| 3-Deep NoRBERT | **3.90** | **22.28** | **16.56** |
| 9-Deep NoRBERT | **3.87** | 19.78 | 13.34 |
| 11-Deep NoRBERT | 3.65 | 18.21 | 11.64 |
| 12-Deep NoRBERT | 3.01 | 16.31 | 9.61 |

Table 3: Masked BLEU score of different models comparing different missing rates.

*5.2. Top NoRBERT*

The above results demonstrate that retraining BERT when we include a GM-VAE layer within may bring imputation improvement when the layer is placed deep inside the BERT network. From this perspective, placing the GMVAE layer in the top of the network, as we do in Top NoRBERT, lacks a priori of any interest. Actually, when we freeze all the parameters from the encoder layers and fine-tune only the classification layer we achieve an imputation accuracy

of 77.14% (Masked) and 75.53% (Disrupted), far below the Deep NoRBERT performance in Table 1. A closer look to the actual imputed words by Top NoRBERT in different sentences led us to conclude that the final GMVAE layer placed right below the classifier promotes topic diversity in the imputation task, which would explain the severe drop in accuracy w.r.t. Deep NoRBERT. This result may be consequence of the fact that upper layers in BERT learn specific features that affect the token choice while the deeper layers pick up general characteristics of text.

Therefore, in order to visualize the effect of our regularizer, Table 4 includes some test sentences reconstructed by Top NoRBERT in comparison with the baseline BERT. In Appendix E we have included more examples with longer sentences (Table E.8) as an extension. For the generation of the results we use again the *snli* dataset with the masking policy defined as 'Low'. The baseline corresponds to BERT model fine-tuned for half an epoch and a learning rate of $5 \cdot 10^{-5}$. The training of Top NoRBERT was fine-tuned with the same configuration. Regarding the GMVAE, we maintained all the previous parameters, except that we increased the learning rate to $10^{-4}$ and trained 200 epochs.

As it is shown in Table 4, the GMVAE stochastic layer at the top of BERT helps it to reconstruct sentences from a robust space, inducing the generation of more diverse sequences than the baseline. It is interesting how it changes some words maintaining the original structure as in the first example in Table 4. Moreover, these alterations maintain grammatical rules ('performs' and 'perform' are used according to the subject) and sometimes correspond to synonymous or analogous words (in this same example, the verb 'sing' is replaced by 'perform', the noun 'choir' by 'band', the object 'masses' by 'friends' and the place 'choir' by 'museum'). This diversity skill is not obtained by the baseline, so it is a characteristic uniquely from our methodology. In other cases, we get changes in words that are not masked so the overall sentence makes sense. The fifth example changes 'out' by 'into' as a consequence of infering 'jumping' from the masked word 'walking'. In the last example, NoRBERT changes 'crossing a overpass' by 'down a intersection sidewalk' as a semantically related structure

17

**Source:** <u>This</u> church <u>choir</u> sings <u>to the masses as they sing joyous songs from the book at a church .</u>

**BERT:** this <span style="color:red">large</span> choir <span style="color:red">looks</span> to the <span style="color:red">camera</span> as they sing <span style="color:red">joy about</span> songs from the book at a church.

**Top-NoRBERT:** <span style="color:red">a dancing band performs</span> to the <span style="color:red">friends</span> as they <span style="color:red">perform funcy bands</span> from the book at a <span style="color:red">museum</span>.

---

**Source:** <u>A man</u> reads <u>the</u> paper <u>in a bar with green lighting</u> .

**BERT:** a man <span style="color:red">in</span> the <span style="color:red">drink</span> in a bar with green lighting.

**Top-NoRBERT:** a man <span style="color:red">on</span> the <span style="color:red">bike</span> in a bar with green <span style="color:red">lights</span>.

---

**Source:** <u>During calf</u> roping <u>a cowboy calls off his</u> horse .

**BERT:** <span style="color:red">a the race</span> a cowboy <span style="color:red">call</span> off his <span style="color:red">back</span>.

**Top-NoRBERT:** during <span style="color:red">horse jumping</span> a cowboy <span style="color:red">tries</span> off his <span style="color:red">dog</span>.

---

**Source:** <u>A man in a black</u> shirt <u>is looking at a</u> bike <u>in a workshop</u> .

**BERT:** a man in a black shirt is looking at a <span style="color:red">woman</span> in a <span style="color:red">conference</span>.

**Top-NoRBERT:** a man in a black shirt is looking at a <span style="color:red">sign</span> in a <span style="color:red">shop</span>.

---

**Source:** <u>The man in the black wetsuit is</u> walking <u>out of the water</u> .

**BERT:** the man in the black wetsuit is <span style="color:red">coming</span> out of the water.

**Top-NoRBERT:** the man in the black <span style="color:red">swimsuit</span> is <span style="color:red">jumping into</span> of the water

---

**Source:** <u>Five girls and two guys are crossing a overpass</u> .

**BERT:** Five girls and two guys are crossing a overpass .

**Top-NoRBERT:** <span style="color:red">three</span> girls and two guys are <span style="color:red">down</span> a <span style="color:red">intersection sidewalk</span>.

---

Table 4: Examples of sentences reconstructed by Top NoRBERT. The first sentence is the original one, with the observed words underlined, i.e. **no underlying means a missing word**. The second is the output of the baseline, BERT fine-tuned. Finally, we show our reconstruction. The words in red correspond to mismatches with the original sentence.

that also corresponds the verb 'to be'.

Enhancing diversity in text generation is a little explored area, as we do not even dispose of clear metrics to measure such an ability, in opposition to for instance image generation, in which researches typically rely on feature space metrics such as the FID to evaluate generation diversity (Heusel et al., 2017). We believe that the Top NoRBERT strategy to achieve such diversity may open future research lines on this topic.

## 6. Extension

Finally, we show how the imputation diversity of traditional seq2seq-type models can be also enhanced by including a regularizer GMVAE layer inside their structure. We start with a simple seq2seq model and then a seq2seq model with attention (Luong et al., 2015).

### 6.1. Seq2seq

A seq2seq model is composed of an encoder which maps the input sentence into a fixed-size vector and a decoder to map this vector into a target sentence. In this architecture, we propose to train a GMVAE over the encoder output as shown in Figure 7. In the fine-tuning step, the encoder is fixed and the decoder is re-trained taking as inputs the GMVAE noisy reconstructed vectors.
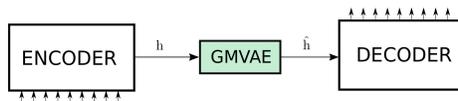


Figure 7: Diagram of the regularized seq2seq model.

### 6.1.1. Results

One of the problems of this first model is caused by the limitations of our baseline. Seq2seq is not suitable for dealing with complex and realistic datasets, that is, long sentences and a wide dictionary, since they encode the semantic and syntactic information of a whole sentence in a single vector, e.g. the encoder

19

output. Notwithstanding, we present in this section some examples where the effect of the regularization layer can be evaluated. The configuration details can be consulted in Appendix C.2.

In Table 5 we show some test sentences reconstructed by our model, compared with the baseline, which is the pre-trained seq2seq model without any GMVAE stochastic layer. In both cases, we reconstruct with the most likely word. By its own, a seq2seq model fulfills its task if the dataset is not very complex, so we have restricted the results to that premise (refer to Appendix B for dataset details). Our model achieves its goal when the sentences are short, finding words that fit the holes, while the baseline fails more often in the task, repeating the previous or following word into the masked place when it does not predict anything better (examples 2, 3 and 5). In addition, our method is able to change other words in the sentence, even if they were not masked, so the overall construction has more sense (example 2). However, in complex scenarios (example 5), both tend to fail, above all our approach, with not grammatically correct sentences.

*6.2. Seq2seq with attention*

We include now the global attention mechanism into the seq2seq (Luong et al., 2015), allowing the network to focus on the relevant parts of the source sentence, acting as an alignment system between encoder and decoder and improving the performance. Now, as the decoder attends to the encoder hidden states at each time step, the previous approach (Section 6.1) results to be insufficient. In this section we present two kind of methodologies: the option 1 regularizes the hidden states in the decoder LSTMs with a Conditional GMVAE (C-GMVAE), and the option 2 the attention vectors with a GMVAE.

The option 1 aims to regularize the hidden states of the decoder at each step $(h_0, h_1, ...h_T)$. To achieve this task, we train a C-GMVAE with pairs of consecutive hidden states $(h_i, h_{i+1})$ from the training sentences, the first one acting as the conditioning input and the second as the input to be reconstructed. See Appendix A for details on the C-GMVAE. At each time step, the C-GMVAE

20

a woman standing in a dark doorway , waiting to be let into the building .

a woman standing in a dark small game waiting to be let into the building .

a woman standing in a dark blue jacket waiting to be let into the building .

a man in an orange hat starring at something .

a man in an hat hat starring at many .

a man in an orange shirt performs at night .

the red car is ahead of the two cars in the background .

the red car is is of the cars cars in the background .

the red car is is of the street cars in the background

five people wearing winter jackets and helmets stand in the snow , with
    snowmobiles in the background.

five girls , winter jackets and helmets stand in the snow , with flowers in
    the background .

five soccer , winter teenager and others stand
    in the snow with this river in the background .

a large bull targets a man , inches away , in a rodeo with his horns , while
    a rodeo clown runs . . .

a bull bull targets a man , petting away , in a bottle with his other , while
    a rodeo clown tries . . .

a young boy move a shoeshine opponent head , wearing a blue with the
    girl , with two boys . . .

Table 5: Examples of sentences reconstructed by the regularized seq2seq. The first sentence is the original one, with the observed words underlined, i.e. **no underlying means a missing word**. The second is the output of the baseline seq2seq pre-trained. Finally, we show our method. The words in red correspond to mismatches with the original sentence.

receives the previous state and the current hidden state $(\hat{h}_{i-1}, h_i)$ to reconstruct the latter $(\hat{h}_i)$. Figure 8a shows a diagram with this approach. We highlight in blue the process concerning the step $i = 1$ as an example, but it is repeated

21

from the beginning until the end-of-sentence token is generated.
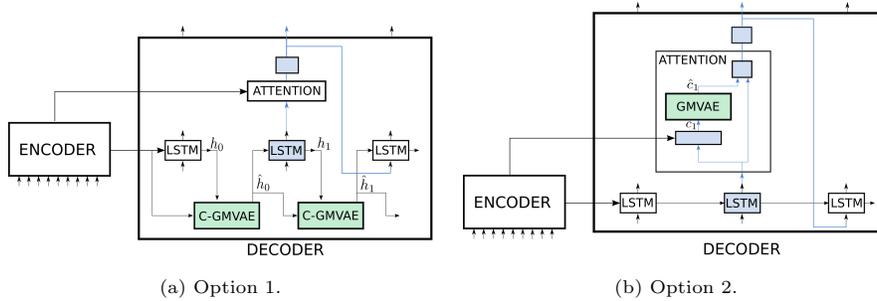


(a) Option 1.  (b) Option 2.

Figure 8: Diagrams of the GMVAE regularized seq2seq model with attention.

Option 2 is structurally simpler. We incorporate the noise in a controlled way, avoiding dependencies on previous states. For that, we propose to introduce the GMVAE layer inside the attention mechanism itself. In particular, the GMVAE layer is trained over the context vectors $(c_0, c_1, ...c_T)$. This model is shown in Figure 8b, where we train the GMVAE with the context vectors of words from the training sentences. We treat each token independently in the GMVAE, since the context vectors usually attend to no more than one or two tokens, thus not requiring a conditional GMVAE.

*6.2.1. Results*

Table 6 shows the results of the two configurations proposed. We use the same dataset as previously but as the model is more powerful due to attention, we are able to increase the percentage of masked tokens to be inferred (see Appendix B for details of this second strategy) without damaging the overall performance of the seq2seq. Moreover, in Appendix E we extend the results for a larger text corpora. The configuration details are presented in Appendix C.3.

The results in Table 6 show how both designs fit our goal, generating new sentences and computing substitutes to the masked tokens that fit the gaps. All of the sentences that are exposed belong to the testing dataset and have been selected randomly. As opposite as in the first scenario in Section 6.1.1, the generation of sentences has improved due to the attention mechanism, so both

22

the baseline and our method perform better the reconstruction of sentences, as was expected. Moreover, the option 2, regularizing the context vectors, not only imputes the masked tokens but also some other tokens in the sentence so the complete structure makes sense. For example, in the second sentence the word 'terrier' is removed and 'dog' is changed of position. More interesting is the third one, where 'kick' and 'stick' are deleted but 'kicking' appears as a verb form of 'kick'.

To understand the diversity of solutions achieved with our model, we can examine not only the most likely imputed word, but also the top five. We focus in option 2 for simplicity. For example, in the first sentence, the baseline best options for 'orange' correspond to colours, however our method also infers the word 'cowboy' in the top 5. In the longest sentence, the forth, we found that even if the final reconstruction was not completely correct (neither in the baseline), our method achieves more varied candidates. In particular, the word 'snowmobiles' has the more likely alternatives ['structure', 'furniture', 'each', 'it' and 'reflections'] for the baseline while ours are ['flags', 'trees', 'umbrellas', 'people' and 'something'], which is a more diverse set that absolutely fits the previous word 'with'.

Our results demonstrate that our proposal performs at least as good as the baseline but in many times is capable to improve generalization in the imputation of missing words. Even more, it can be seen as a way of data augmentation in the sense that builds new sentences, acceptable and different from the baseline choices.

## 7. Conclusions and future work

In this work we have proved the successful effect of adding a stochastic GMVAE layer in BERT through NoRBERT. We study the different advantages regarding the layer where it is applied. While Top NoRBERT successes with an increment of diversity as well as an easier way of adaptability to new contexts, Deep NoRBERT responds better in terms of accuracy and BLEU score. In the

former case, we propose a novel methodology to generate new structures of text with diverse topics that fit the gaps thanks to the inclusion of controlled noise through a DGM. As a way of reinforcing our idea, we prove the GMVAE effect regularizing a well-studied scenario with FMNIST images.

As an extension, in Section 6 we present the advantages of the stochastic layer in autoregressive seq2seq models. Despite their limitations with long sentences, our method is able to predict assorted structures upon an extend. Then, we enforced the same idea applying attention and exploring other scenarios that incorporate the regularization at different points of the baseline. In this work, we successfully reconstruct a varied set of topics from the masked source sentences and demonstrate the efficacy of the stochastic layer in finding synonymous or analogous fragments that fit in the gaps.

For now, there is no metric to evaluate robust and varied solutions, since traditional evaluations as BLEU (Papineni et al., 2002) or ROUGE (Lin and Och, 2004) are based in the reconstruction of the original sentence. There is no perfect evaluation metrics testing the text generation because it is difficult to resume all the semantic and syntactic properties that language needs to fulfil (Wang et al., 2019b). Therefore, we let for future work the exploration of metrics or lost functions that allows the LM to generate sentence embeddings with more diversity based on the context.

**Acknowledgements**

# References

Bahdanau D, Cho K, Bengio Y. Neural machine translation by jointly learning to align and translate. In: 3rd International Conference on Learning Representations, ICLR 2015. 2015. .

[dataset] Bowman S, Angeli G, Potts C, Manning CD. A large annotated corpus for learning natural language inference. In: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. 2015. p. 632–42.

Bowman S, Vilnis L, Vinyals O, Dai A, Jozefowicz R, Bengio S. Generating sentences from a continuous space. In: Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning. 2016. p. 10–21.

Caccia M, Caccia L, Fedus W, Larochelle H, Pineau J, Charlin L. Language gans falling short. In: Proceedings of the Eighth International Conference on Learning Representations, ICLR 2020. 2020. .

Dai Z, Yang Z, Yang Y, Carbonell JG, Le Q, Salakhutdinov R. Transformer-xl: Attentive language models beyond a fixed-length context. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. 2019. p. 2978–88.

Devlin J, Chang MW, Lee K, Toutanova K. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:181004805 2018;.

Dilokthanakul N, Mediano PA, Garnelo M, Lee MC, Salimbeni H, Arulkumaran K, Shanahan M. Deep unsupervised clustering with gaussian mixture variational autoencoders. arXiv preprint arXiv:161102648 2016;.

Donahue D, Rumshisky A. Adversarial text generation without reinforcement learning. arXiv preprint arXiv:181006640 2018;.

Dong L, Yang N, Wang W, Wei F, Liu X, Wang Y, Gao J, Zhou M, Hon HW. Unified language model pre-training for natural language understanding and

generation. In: Advances in Neural Information Processing Systems. 2019. p. 13042–54.

[dataset] Elliott D, Frank S, Sima'an K, Specia L. Multi30k: Multilingual english-german image descriptions. In: Proceedings of the 5th Workshop on Vision and Language. 2016. p. 70–4.

Fedus W, Goodfellow I, Dai AM. Maskgan: Better text generation via filling in the _. In: International Conference on Learning Representations. 2018. .

Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville A, Bengio Y. Generative adversarial nets. In: Advances in neural information processing systems. 2014. p. 2672–80.

Gu J, Wang C, Zhao J. Levenshtein transformer. In: Advances in Neural Information Processing Systems. 2019. p. 11179–89.

Gulcehre C, Firat O, Xu K, Cho K, Barrault L, Lin HC, Bougares F, Schwenk H, Bengio Y. On using monolingual corpora in neural machine translation. arXiv preprint arXiv:150303535 2015;.

Guo J, Lu S, Cai H, Zhang W, Yu Y, Wang J. Long text generation via adversarial training with leaked information. In: Thirty-Second AAAI Conference on Artificial Intelligence. 2018. .

Gupta A, Agarwal A, Singh P, Rai P. A deep generative framework for paraphrase generation. In: Thirty-Second AAAI Conference on Artificial Intelligence. 2018. .

Haidar MA, Rezagholizadeh M. Textkd-gan: Text generation using knowledge distillation and generative adversarial networks. In: Canadian Conference on Artificial Intelligence. Springer; 2019. p. 107–18.

Haidar MA, Rezagholizadeh M, Do Omri A, Rashid A. Latent code and text-based generative adversarial networks for soft-text generation. In: Proceedings of the 2019 Conference of the North American Chapter of the Associa-

tion for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers). 2019. p. 2248–58.

Heusel M, Ramsauer H, Unterthiner T, Nessler B, Hochreiter S. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In: Advances in neural information processing systems. 2017. p. 6626–37.

Hu Z, Yang Z, Liang X, Salakhutdinov R, Xing EP. Toward controlled generation of text. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR. org; 2017. p. 1587–96.

Jang E, Gu S, Poole B. Categorical reparametrization with gumble-softmax. In: International Conference on Learning Representations (ICLR 2017). OpenReview. net; 2017. .

Jiang Z, Zheng Y, Tan H, Tang B, Zhou H. Variational deep embedding: an unsupervised and generative approach to clustering. In: Proceedings of the 26th International Joint Conference on Artificial Intelligence. 2017. p. 1965–72.

Keskar NS, McCann B, Varshney LR, Xiong C, Socher R. Ctrl: A conditional transformer language model for controllable generation. arXiv preprint arXiv:190905858 2019;.

Kingma DP, Welling M. Auto-encoding variational bayes. stat 2014;1050:1.

Kovaleva O, Romanov A, Rogers A, Rumshisky A, Romanov A, De-Arteaga M, Wallach H, Chayes J, Borgs C, Chouldechova A, et al. Revealing the dark secrets of bert. In: Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). American Medical Informatics Association; volume 1; 2019. p. 2465–75.

Krogh A, Hertz JA. A simple weight decay can improve generalization. In: Advances in neural information processing systems. 1992. p. 950–7.

Lample G, Conneau A. Cross-lingual language model pretraining. arXiv preprint arXiv:190107291 2019;.

Lin CY, Och F. Looking for a few good metrics: Rouge and its evaluation. In: Ntcir Workshop. 2004. .

Liu Y, Ott M, Goyal N, Du J, Joshi M, Chen D, Levy O, Lewis M, Zettlemoyer L, Stoyanov V. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:190711692 2019;.

Loper E, Bird S. Nltk: The natural language toolkit. In: Proceedings of the ACL-02 Workshop on Effective Tools and Methodologies for Teaching Natural Language Processing and Computational Linguistics. 2002. p. 63–70.

Lu S, Zhu Y, Zhang W, Wang J, Yu Y. Neural text generation: Past, present and beyond. arXiv preprint arXiv:180307133 2018;.

Luong MT, Pham H, Manning CD. Effective approaches to attention-based neural machine translation. In: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. 2015. p. 1412–21.

Ma X, Zhang P, Zhang S, Duan N, Hou Y, Zhou M, Song D. A tensorized transformer for language modeling. In: Advances in Neural Information Processing Systems. 2019. p. 2229–39.

de Masson d'Autume C, Mohamed S, Rosca M, Rae J. Training language gans from scratch. In: Advances in Neural Information Processing Systems. 2019. p. 4302–13.

Pagnoni A, Liu K, Li S. Conditional variational autoencoder for neural machine translation. arXiv preprint arXiv:181204405 2018;.

Papineni K, Roukos S, Ward T, Zhu WJ. Bleu: a method for automatic evaluation of machine translation. In: Proceedings of the 40th annual meeting on association for computational linguistics. Association for Computational Linguistics; 2002. p. 311–8.

Radford A, Metz L, Chintala S. Unsupervised representation learning with deep convolutional generative adversarial networks. arXiv preprint arXiv:151106434 2015;.

Radford A, Narasimhan K, Salimans T, Sutskever I. Improving language understanding by generative pre-training. unpublished work 2018;URL: https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/lang

Radford A, Wu J, Child R, Luan D, Amodei D, Sutskever I. Language models are unsupervised multitask learners. unpublished work 2019;URL: https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/language-unsupervised/la

Rashid A, Do-Omri A, Haidar MA, Liu Q, Rezagholizadeh M. Bilingual-gan: A step towards parallel text generation. In: Proceedings of the Workshop on Methods for Optimizing and Evaluating Neural Language Generation. 2019. p. 55–64.

Sanh V, Debut L, Chaumond J, Wolf T. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. arXiv preprint arXiv:191001108 2019;.

Shen D, Zhang Y, Henao R, Su Q, Carin L. Deconvolutional latent-variable model for text sequence matching. In: Thirty-Second AAAI Conference on Artificial Intelligence. 2018. .

Shi W, Zhou H, Miao N, Zhao S, Li L. Fixing gaussian mixture vaes for interpretable text generation. arXiv preprint arXiv:190606719 2019;.

Sohn K, Lee H, Yan X. Learning structured output representation using deep conditional generative models. In: Advances in neural information processing systems. 2015. p. 3483–91.

Song K, Tan X, Qin T, Lu J, Liu TY. Mass: Masked sequence to sequence pre-training for language generation. In: International Conference on Machine Learning. 2019. p. 5926–36.

Sriram A, Jun H, Satheesh S, Coates A. Cold fusion: Training seq2seq models together with language models. Proc Interspeech 2018 2018;:387–91.

Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov R. Dropout: a simple way to prevent neural networks from overfitting. The journal of machine learning research 2014;15(1):1929–58.

Subramanian S, Mudumba SR, Sordoni A, Trischler A, Courville AC, Pal C. Towards text generation with adversarially learned neural outlines. In: Advances in Neural Information Processing Systems. 2018. p. 7551–63.

Sun Y, Wang S, Li Y, Feng S, Tian H, Wu H, Wang H. Ernie 2.0: A continual pre-training framework for language understanding. In: Thirty-Fourth AAAI Conference on Artificial Intelligence. 2020. .

Sutskever I, Vinyals O, Le QV. Sequence to sequence learning with neural networks. In: Advances in neural information processing systems. 2014. p. 3104–12.

Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser Ł, Polosukhin I. Attention is all you need. In: Advances in neural information processing systems. 2017. p. 5998–6008.

Wang A, Singh A, Michael J, Hill F, Levy O, Bowman S. Glue: A multi-task benchmark and analysis platform for natural language understanding. In: 7th International Conference on Learning Representations, ICLR 2019. 2019a. .

Wang B, Wang A, Chen F, Wang Y, Kuo CCJ. Evaluating word embedding models: methods and experimental results. APSIPA Transactions on Signal and Information Processing 2019b;8.

Wu Y, Schuster M, Chen Z, Le QV, Norouzi M, Macherey W, Krikun M, Cao Y, Gao Q, Macherey K, et al. Google's neural machine translation system: Bridging the gap between human and machine translation. arXiv preprint arXiv:160908144 2016;.

Xiao H. Hungarian layer: A novel interpretable neural layer for paraphrase identification. Neural Networks 2020;131:172–84.

[dataset] Xiao H, Rasul K, Vollgraf R. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. arXiv preprint arXiv:170807747 2017;.

Yang M, Chen L, Lyu Z, Liu J, Shen Y, Wu Q. Hierarchical fusion of common sense knowledge and classifier decisions for answer selection in community question answering. Neural Networks 2020a;132:53–65.

Yang S, Lu H, Kang S, Xue L, Xiao J, Su D, Xie L, Yu D. On the localness modeling for the self-attention based end-to-end speech synthesis. Neural Networks 2020b;125:121–30.

Yang Z, Dai Z, Yang Y, Carbonell J, Salakhutdinov RR, Le QV. Xlnet: Generalized autoregressive pretraining for language understanding. In: Advances in neural information processing systems. 2019. p. 5754–64.

Yang Z, Hu Z, Salakhutdinov R, Berg-Kirkpatrick T. Improved variational autoencoders for text modeling using dilated convolutions. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR. org; 2017. p. 3881–90.

Yu L, Zhang W, Wang J, Yu Y. Seqgan: sequence generative adversarial nets with policy gradient. In: Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence. 2017. p. 2852–8.

Yu W, Zheng C, Cheng W, Aggarwal CC, Song D, Zong B, Chen H, Wang W. Learning deep network representations with adversarially regularized autoencoders. In: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2018. p. 2663–71.

Zhang H, Goodfellow I, Metaxas D, Odena A. Self-attention generative adversarial networks. In: International Conference on Machine Learning. PMLR; 2019. p. 7354–63.

Zhang Y, Gan Z, Fan K, Chen Z, Henao R, Shen D, Carin L. Adversarial feature matching for text generation. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR. org; 2017. p. 4006–15.

Zhao J, Kim Y, Zhang K, Rush AM, LeCun Y. Adversarially regularized autoencoders. In: 35th International Conference on Machine Learning, ICML 2018. International Machine Learning Society (IMLS); 2018. p. 9405–20.

Zhu Y, Kiros R, Zemel R, Salakhutdinov R, Urtasun R, Torralba A, Fidler S. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In: Proceedings of the IEEE international conference on computer vision. 2015. p. 19–27.

## Appendix A. The Conditional GMVAE

When we deal with conditional DGM, we mean that the entire generative process is conditioned on some extra observed inputs. Sohn et al. (2015) presented Conditional Variational Autoencoder (CVAE), where the observations modulate the Gaussian prior. In a similar way, we have studied two architectures to condition our distributions on an input that we have defined $h$. In this section, we expose the changes applied and describe the two versions of C-GMVAE that we have explored, referred as models A and B.



(a) C-GMVAE model A.      (b) C-GMVAE model B.

Figure A.9: The directed graphical models considered for the C-GMVAE in the work. Solid lines denote the generative model and dashed lines the variational approximation.

The first architecture (model A) that we tried is shown in Figure A.9a. For its implementation we had to change the prior distribution of $w$ as $p(w|h) \sim \mathcal{N}(\mu(h), \Sigma(h))$, where the mean and variance of the normal distribution are parameterized by dense nets, and $q_{\phi_w}(w|x, h) = \mathcal{N}(\mu_{\phi_w}([x; h]), \Sigma_{\phi_w}([x; h]))$, where we only concatenate $h$ to the original input $x$. The main drawback of this model is that the reconstruction as we performed it (compute $z$ from $x$ through the inference model and then reconstruct $x$ from this $z$ by the generative model) does not use the observed $h$.

The graph in Figure A.9b belongs to our second version (model B), the one we applied to the presented results. In contrast, in the generative model, we now maintain the original prior of $w$ but condition the $z$ distribution on $h$

as Equation A.1a. For the variational family, we modify the encoding of $z$ as Equation A.1b.

$$p_{\beta_z}(z|w,y,h) = \prod_{k=1}^{K} \mathcal{N}(\mu_{\beta_z}([w;h]), \Sigma_{\beta_z}([w;h]))^{y_k==1} \tag{A.1a}$$

$$q_{\phi_z}(z|x,h) = \prod_{k=1}^{K} \mathcal{N}(\mu_{\phi_z}([x;h]), \Sigma_{\phi_z}([x;h])) \tag{A.1b}$$

## Appendix B. Dataset for the extension results

In Sections 6.1.1 and 6.2.1 we train the models with the **multi30k** dataset ([dataset] Elliott et al., 2016). It consists on a training set of 29000 English sentences and a set of 1000 test sentences. The vocabulary size is 10118 tokens. It is not a very large corpora, but we mask some tokens so the scenario gets more complicated to be trained. We use different strategies in the masking process, with higher and lower rates.

In the first strategy, we use a policy of masked tokens more sophisticated that permit the masking of a less percentage of words but focusing on nouns, verbs, adjectives... That is, ignoring stopwords. We use the English stopwords list from the *nltk*[6] library (Loper and Bird, 2002). We mask the 80% of the sentences and we generate two masks in a sentence with a probability of 0.8. Among these, we also generate a third mask with probability of 0.8.

In the second strategy, we increase the number of masked tokens and do not exclude any type of grammatical word so any one can be deleted. In this policy, we mask each token with a probability of 0.6. Therefore, we have more [MASK] tokens than proper words.

---

[6]https://www.nltk.org/

## Appendix  C. Configuration and experiments

*Appendix  C.1. FMNIST*

The model used for the experiments with the FMNIST dataset consists on 9 linear layers with RELU as the activation function. The size of the output features on each layer is, from bottom to top, 700, 600, 512, 256, 128, 64, 32, 16 and 10, which corresponds with the number of classes. We employ a negative log-likelihood loss function and the stochastic gradient descent with a learning rate of 0.01 for the optimization.

The dataset is composed of 28x28 images in grayscale associated with a label from 10 classes. We divide the set in 12000 samples for training and 48000 for validation. The test set has 10000 images. The only preprocessing step is the normalization to 0.5 mean and variance.

*Appendix  C.2. Seq2seq*

For this model, we follow the networks structure from Luong et al. (2015), omitting the attention mechanism for now. Consequently, we will use a LSTM as the RNN unit, a bidirectional encoder, a depth of two layers in the networks and a hidden size of 1024 for each of them.

The configuration of this model follows a seq2seq pre-training of 120 epochs and a fine-tuning of the regularized decoder for only 20 epochs after training the GMVAE. In the GMVAE, after different proves we finally chose 1500 for the hidden dimension, 100 for $z$, 20 for $w$ and a $K$ of 10 MoG of the prior. The depth in the networks is 5 layers and the deviation, $\sigma$, of the posterior normal distribution in the decoder $10^{-4}$. We saved the hidden states (encoder output) of all the training sentences, and trained the GMVAE for 100 epochs.

*Appendix  C.3. Seq2seq with attention*

We keep the same configuration from Luong et al. (2015), but including the global attention mechanism.

In the option 1, the C-GMVAE is trained with consecutive pairs of hidden states from the whole training set during 100 epochs. We finally used a hidden

dimension of 1500, 150 for the latent space of $z$ and 50 for $w$. We configured $K = 20$ classes in the MoG and a $\sigma$ of $10^{-2}$ for the decoder posterior. The number of layers on each of the distributions modeled was 6. During the training we selected a learning rate of $10^{-5}$, a dropout of 0.3 and a batch size of 64.

In the option 2, the GMVAE is trained for 150 epochs with the same configuration as before. It only changes the graph as described in Section 2.1.

In both options, for the pre-training of the seq2seq, 30 epochs were enough since the attention mechanism eases the convergence of the model. After the training of the C-GMVAE and the GMVAE respectively, we fine-tuned the seq2seq decoder with the inclusion of the suitable stochastic layer as mentioned in Section 6.2 during other 30 epochs.


## Appendix D. Other models

*Appendix D.1. Seq2seq with attention*

In the autoregressive model of seq2seq with attention, we, firstly, tried training the C-GMVAE from Figure A.9a to generate each hidden state conditioned on the previous one. These generated states were the inputs for the next LSTM unit. However, it did not work as good as we expected. Consequently, we changed the process in a way that instead of using the DGM to generate samples, we could take advantage of its latent space and reconstruct the original hidden states from the LSTMs.

Inspired by the first model (Section 6.1), we also tried to follow the same idea that is presented as option 1 in Section 6.2 but conditioning always in the encoder output instead of the previous hidden state, but it did not improve neither the results that we are presenting in this work.

Regarding the option 2, initially we used a simpler approach, regularizing the attention output after concatenating it with the LSTM output and exactly before applying the classification layer that matches the vocabulary size. Here, the training was not successful and the imputed words did not follow grammatical rules as a LM is expected to do. After this, we tried the integration of the

GMVAE in a previous step, as it is successfully explained in this work.

## Appendix E. Additional results

Table E.7 presents additional results from the option 2 in the seq2seq with attention model using the *snli* dataset. Once again, we prove the efficacy of our method, even if the dataset gets more complicated. In this table we present different samples of sentences reconstructed from the masked template, following the same philosophy of the results in Table 6. The fifth example exposes an extreme case where it is only observed the first word, 'a', and both the baseline and our method infer completely different sequences but good alternatives at the same time.

Table E.8 is an extension of Table 4 with results form Top NoRBERT.

a <u>man</u> in an orange <u>hat starring at</u> something .

a <u>man</u> in a <span style="color:red">black</span> <u>hat starring at</u> something .

a <u>man</u> in a <span style="color:red">hard</span> <u>hat starring at</u> something .

a <u>man</u> in a <span style="color:red">black</span> <u>hat starring at</u> something .

---

a boston <u>terrier</u> is <u>running</u> on lush <u>green grass</u> in <u>front of</u> a white fence .

a <span style="color:red">gray</span> <u>terrier</u> <span style="color:red">dog</span> <u>running</u> on <span style="color:red">the</span> <u>green grass</u> in <u>front of</u> a <span style="color:red">blue shack</span> .

a <span style="color:red">gray</span> <u>terrier</u> <span style="color:red">dog</span> <u>running</u> <span style="color:red">through tall</span> <u>green grass</u> in <u>front of</u> a <span style="color:red">red ball</span>

a <span style="color:red">black dog</span> is <u>running</u> <span style="color:red">through the grass</span> <u>grass</u> in <u>front of</u> a <span style="color:red">red flag</span> .

---

a girl in karate <u>uniform</u> breaking <u>a stick</u> with a front <u>kick .</u>

a <span style="color:red">man in a</span> <u>uniform</u> <span style="color:red">throws</span> <u>a stick</u> <span style="color:red">to his his</span> <u>kick .</u>

a <span style="color:red">boy in a</span> <u>uniform</u> <span style="color:red">with</span> <u>a stick</u> <span style="color:red">in a large</span> <u>kick .</u>

a <span style="color:red">man in a</span> <u>uniform</u> <span style="color:red">kicking</span> <u>a</u> <span style="color:red">ball up to his opponent</span> .

---

five people wearing <u>winter</u> jackets and helmets <u>stand</u> in the snow , <u>with</u> snowmobiles in the <u>background</u> .

<span style="color:red">two men in</span> <u>winter</u> jackets and <span style="color:red">hats</span> stand in <span style="color:red">a large space</span> with <span style="color:red">structure</span> in the background .

<span style="color:red">a group of</span> <u>winter</u> <span style="color:red">day at a</span> stand in <span style="color:red">a snowy area</span> with <span style="color:red">trees</span> in the background

<span style="color:red">two men wearing</span> <u>winter</u> <span style="color:red">clothing</span> and <span style="color:red">hats</span> stand <span style="color:red">on the snow covered street</span> with <span style="color:red">flags open</span> .

---

<u>a</u> man in a <u>vest</u> is <u>sitting</u> in a chair <u>and</u> holding magazines .

<u>a</u> man in a <u>vest</u> is <u>sitting</u> <span style="color:red">on</span> a <span style="color:red">rock</span> <u>and</u> <span style="color:red">looking out</span> .

<u>a</u> man in a <u>vest</u> is <u>sitting</u> <span style="color:red">on</span> a <span style="color:red">sidewalk</span> <u>and</u> <span style="color:red">playing music</span> .

<u>a</u> man <span style="color:red">wearing</span> a <u>vest</u> is <u>sitting</u> <span style="color:red">on</span> a <span style="color:red">wall</span> <u>and</u> <span style="color:red">smoking a cigarette</span> .

---

a <u>mother</u> and <u>her</u> young son <u>enjoying</u> a beautiful <u>day outside</u> .

a <u>mother</u> and <u>her</u> <span style="color:red">daughter are</span> <u>enjoying</u> a <span style="color:red">wedding</span> day outside .

a <u>mother</u> and <u>her</u> <span style="color:red">child are</span> <u>enjoying</u> a <span style="color:red">hot</span> day outside .

a <u>mother</u> and <u>her</u> <span style="color:red">children are</span> <u>enjoying</u> a <span style="color:red">hot</span> day outside .

---

Table 6: Examples of sentences reconstructed by the regularized seq2seq with attention following the same format as Table 5: original, baseline, options 1 and 2.

an old man with a package poses in front of an advertisement .

an old man is standing with arms in front of an audience .

an old man in a blue shirt in front of an audience .

---

a man playing an electric guitar on stage .

a man playing an electric guitar on stage .

a man plays an electric guitar and sings .

---

a blond-haired doctor and her african american assistant looking threw
    new medical manuals .

a man is standing in an american assistant , using a medical apparatus .

a man is looking at the american nurse to get a medical patient .

---

a young family enjoys feeling ocean waves lap at their feet .

a young boy is feeling ocean and is on the beach .

a young man in feeling ocean is surfing on a surfboard .

---

a man reads the paper in a bar with green lighting .

a man is standing in front of a crowd of people .

a man is sitting on a bench reading a book while sitting

---

three firefighter come out of subway station .

three people come down a street corner .

three people come out of a boat .

---

a person wearing a straw hat , standing outside working a steel apparatus
    with a pile of coconuts on the ground .

a man wearing a straw hat , standing outside of a steel structure with a
    blue umbrella laying on the ground .

a man wearing a straw hat , standing outside a large steel structure with a
    tree in front of the ground .

---

Table E.7: Additional examples of sentences reconstructed by the regularized hidden states in the seq2seq with attention. Sentences order: original, baseline and reconstruction from our regularized option 2 of the GMVAE and the context vectors.

39

A man looking over a bicycle 's rear wheel in the maintenance garage with various tools visible in the background .

a man looking over a bicycle's back wheel in the maintenance garage with his tools visible in the background.

a man looking over a bicycle's rear wheel in a construction garage with wooden equipment is in the background.

---

A person dressed in a dress with flowers and a stuffed bee attached to it , is pushing a baby stroller down the street .

a person dressed in a suit with flowers and a stuffed animal attached to it, is pushing a baby stroller down the street.

a person dressed in a shirt with flowers and pink stuffed toy over to it, is riding a baby stroller down the street.

---

A blond-haired doctor and her African american assistant looking threw new medical manuals .

a blond - haired doctor and her african american doctor looking at new medical scrubs.

a blond - haired nurse and her african asian owner looking around new medical equipments.

---

3 young man in hoods standing in the middle of a quiet street facing the camera .

3 young man in hoods standing in the middle of a busy street facing the camera.

a young man in sunglassess standing in the front of a busy street holding the camera.

---

Table E.8: Additional examples of sentences reconstructed by Top NoRBERT. The first sentence is the original one, with the observed words underlined. The second is the output of the baseline BERT fine-tuned. Finally, we show the reconstruction. The words in red correspond to mismatches with the original sentence.