

# FINE-TUNING PRETRAINED MODELS WITH NVIB FOR IMPROVED GENERALISATION

Fabio Fehr<sup>1,2</sup>    Alina Elena Baia<sup>\*1</sup>    Xiaoguang Chang<sup>\*†1,3</sup>    Andrei C. Coman<sup>\*1,2</sup>  
 Karl El Hajal<sup>\*1,2</sup>    Dina El Zein<sup>\*1,2</sup>    Shashi Kumar<sup>\*1,2</sup>    Juan Zuluaga-Gomez<sup>\*1,2</sup>  
 Andrea Cavallaro<sup>1,2</sup>    Damien Teney<sup>1</sup>    James Henderson<sup>1</sup>

<sup>1</sup>Idiap Research Institute, Switzerland

<sup>2</sup>Ecole Polytechnique Fédérale de Lausanne, Switzerland

<sup>3</sup>Southeast University, China

first.last@idiap.ch, †xg\_chang@seu.edu.cn

## ABSTRACT

Fine-tuned pretrained attention-based models often struggle with generalisation, leading to poor performance on tasks like out-of-domain transfer, distribution shifts, and few-shot learning. This limitation is prevalent across modalities such as speech, text, graphs, and vision. Nonparametric Variational Information Bottleneck (NVIB) is an attention-based information-theoretic regulariser applicable to pretrained models that has been shown to improve generalisation. However, prior work has applied NVIB only to the text modality and without fine-tuning. We investigate whether NVIB’s ability to remove information from pretrained embeddings helps the model avoid spurious correlations with noisy and superficial features during fine-tuning. We are the first to integrate NVIB regularisation during fine-tuning across multiple diverse models and modalities. This required modifications to the architecture which enhance adaptability and stability during fine-tuning and simplify the evaluation. We found improved out-of-distribution generalisation in: speech quality assessment and language identification, text with induced attention sparsity, graph-based link prediction, and few-shot image classification. <sup>1</sup>

## 1 INTRODUCTION

Leveraging pretrained attention-based representations by fine-tuning has become the de facto modelling paradigm due to its wide applicability and significant improvements on the state-of-the-art (Ruder et al., 2019). Applications of pretrained Transformers (Vaswani et al., 2017) are modality agnostic and gained prevalence across: speech processing (Baevski et al., 2020; Radford et al., 2023); natural language processing (Devlin et al., 2019; Raffel et al., 2020; Touvron et al., 2023), graphs Rong et al. (2020); Li et al. (2021b) and computer vision (Liu et al., 2021; Dosovitskiy et al., 2021; Bao et al., 2022).

The success of pretrained attention-based models is thought to stem from their ability to scale, both in terms of corpora size and the number of parameters, as well as the inductive biases inherent in the attention-based architecture (Henderson, 2020; Zhai et al., 2022; Fedus et al., 2021; Dehghani et al., 2023). Despite their success, these models still exhibit notable limitations during fine-tuning. Due to their large number of parameters and expressivity, they can be prone to overfitting and struggle to generalise in the presence of shortcuts from spurious correlations (Bhargava et al., 2021; Geirhos et al., 2020) and distribution shift (Wu et al., 2020a; Kumar et al., 2022). The attention mechanism

<sup>1</sup>The code is publically available at:

<https://github.com/idiap/nvib> &  
[https://github.com/idiap/nvib\\_finetuning](https://github.com/idiap/nvib_finetuning)

\*Equal contribution, alphabetical order.

facilitates expressivity through token interaction, but this also introduces redundant information, which can hinder generalisation (Bian et al., 2021; Bhojanapalli et al., 2021). Introducing sparsity as a form of regularisation into attention has been shown to improve generalisation performance by reducing this redundancy (Child et al., 2019; Behjati et al., 2023; Fehr & Henderson, 2024). However, regularising attention during fine-tuning of pretrained models remains both challenging and unexplored.

Information bottleneck (IB) is an information-theoretic regulariser that learns latent features  $Z$  that compress the input  $X$  while preserving information for the downstream task  $Y$  (Tishby et al., 2000). The variational information bottleneck (VIB) framework, introduced through a variational lower bound to the IB objective (Alemi et al., 2017), enables deep neural representations (Tishby & Zaslavsky, 2015) to be trained using gradient-based optimisation. This framework has been widely applied across speech (Nelus & Martin, 2021; Lian et al., 2022), natural language (McCarthy et al., 2020; mahabadi et al., 2021), graphs (Wu et al., 2020b; Sun et al., 2022) and vision (Han et al., 2020; Chun, 2024). The success of the VIB framework can be attributed to its key properties, including resilience against spurious correlations (Chuah et al., 2022) and distribution shift (Li et al., 2021a), robustness (Zhang et al., 2022) and sparsity (Paranjape et al., 2020). Despite this success, VIB regularisation has seen limited exploration in the fine-tuning of pretrained attention-based models. Applying VIB to these pretrained models is difficult due to the complexity of incorporating it into the variable-sized latent representations accessed by attention.

Henderson & Fehr (2023) propose Nonparametric Variational Information Bottleneck (NVIB) as a VIB regulariser for attention layers. NVIB regularises the variable-sized representations accessed by attention by compressing both the information in individual vectors and the number of vectors. Further contributions to NVIB have demonstrated characteristics such as out-of-distribution (OOD) generalisation and sparsity (Henderson & Fehr, 2023; Behjati et al., 2023; Fehr & Henderson, 2024). Behjati et al. (2023) employ NVIB for representation learning by incorporating the regulariser into the self-attention layers of a Transformer-based encoder, and trains from scratch to progressively learn sparser representations through its layers. Fehr & Henderson (2024) integrated NVIB into pretrained models and achieved improvements in OOD summarisation and translation tasks without further training. Previous work has not applied NVIB regularisation during fine-tuning of pretrained models, nor has it explored generalising nonparametric variational models beyond text to diverse modalities like vision, speech, and graphs with their varying model architectures, data, and tasks.

**Contributions.** In this paper, we are the first to extend NVIB regularisation methods to fine-tuning, with diverse pretrained models. (1) We propose several novel methods for NVIB fine-tuning, including a learnable prior mean embedding per layer for adaptability, clipped Dirichlet pseudo-counts for stability, and a simplified denoising attention function at evaluation (Section 2). (2) We do the first empirical evaluation of NVIB on diverse modalities such as speech (Section 3.1), text (Section 3.2), graphs (Section 3.3), and vision (Section 3.4). (3) We show improved OOD generalisation in classification and regression tasks, demonstrating NVIB’s added value across diverse applications.

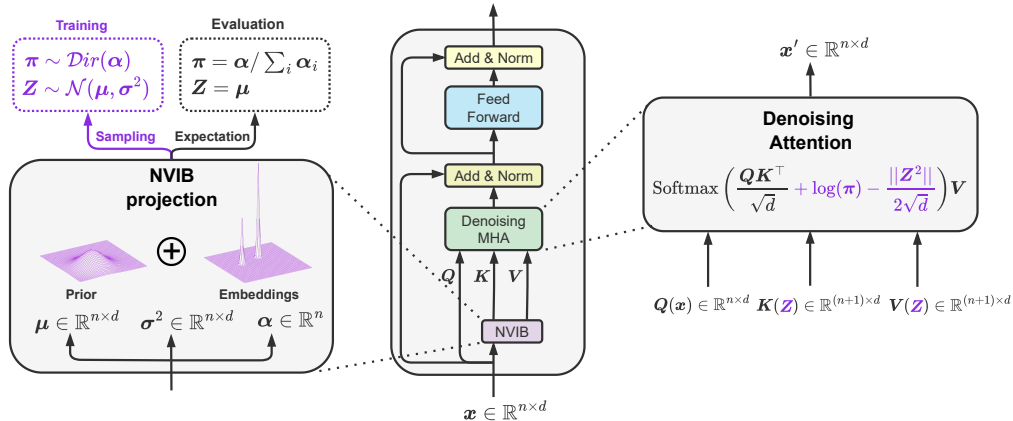


Figure 1: The NVIB module including the NVIB layer (left) and denoising attention (right).

## 2 FINE-TUNING WITH NVIB

Figure 1 depicts an NVIB module, with the NVIB layer (left) and denoising attention function (right). The NVIB layer projects the sequence of vectors  $\mathbf{x} \in \mathbb{R}^{n \times d}$  from a Transformer embedding to the parameters of a Dirichlet Process. These parameters include the isotropic Gaussian means  $\boldsymbol{\mu} \in \mathbb{R}^{(n+1) \times d}$  and variances  $\boldsymbol{\sigma}^2 \in \mathbb{R}^{(n+1) \times d}$ , and the Dirichlet concentration parameters  $\boldsymbol{\alpha} \in \mathbb{R}^{(n+1)}$ . Each of the  $n$  vectors has an associated mixture component, along with an additional  $(n+1)^{\text{th}}$  component that serves as a prior for the embeddings. During training, the NVIB layer samples a mixture distribution, represented as a set of weighted vectors  $(\boldsymbol{\pi}, \mathbf{Z})$ , where  $\boldsymbol{\pi} \sim \text{Dir}(\boldsymbol{\alpha})$  and  $\mathbf{Z}_i \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\sigma}_i^2)$ . During evaluation, the NVIB layer outputs the expectation of these samples, which is the mixture of  $n+1$  Gaussians, but can be approximated as  $\mathbf{Z} = \boldsymbol{\mu}$  and  $\boldsymbol{\pi} = \boldsymbol{\alpha} / \sum_i \alpha_i$ .

Figure 1 (right) depicts how the denoising attention function is a generalisation of standard attention to any nonparametric mixture distribution. In the case of a set of weighted vectors, this involves using the weights  $\boldsymbol{\pi}$  as bias terms for the attention weights over keys  $K(\mathbf{Z})$ . We provide a detailed description and pseudocode for denoising attention in Appendix B, and a consolidated overview of prior research on NVIB in Appendix A.

Following from Fehr & Henderson (2024), we reinterpret the pretrained models as nonparametric variational models by including NVIB layers before the attention mechanisms. This layer maps the input vectors  $\mathbf{x}$  to the DP parameters  $(\boldsymbol{\mu}^q, \boldsymbol{\sigma}^q, \boldsymbol{\alpha}^q)$ :

$$\begin{aligned} \boldsymbol{\mu} &= \boldsymbol{\mu}(\mathbf{x}) = \mathbf{x}\mathbf{W}^\mu + \mathbf{b}^\mu & \mathbf{W}^\mu &= \mathbf{I}; & \mathbf{b}^\mu &= \mathbf{0} \\ \boldsymbol{\sigma}^2 &= \boldsymbol{\sigma}^2(\mathbf{x}) = \exp(\mathbf{x}\mathbf{W}^\sigma + \mathbf{b}^\sigma) & \mathbf{W}^\sigma &= \mathbf{0}; & \mathbf{b}^\sigma &= \log(\tau_\sigma^2) \\ \boldsymbol{\alpha} &= \boldsymbol{\alpha}(\mathbf{x}) = \exp(\mathbf{x}^2\mathbf{w}_1^\alpha + \mathbf{x}\mathbf{w}_2^\alpha + b^\alpha) & \mathbf{w}_1^\alpha &= \frac{1}{2\sqrt{d/h}} \odot \mathbf{1}; & \mathbf{w}_2^\alpha &= \mathbf{0}; & b^\alpha &= \tau_\alpha \end{aligned}$$

This initialisation ensures empirical equivalence with the pretrained model, after manual adjustment of the hyperparameters  $(\tau_\sigma^2, \tau_\alpha)$  for each model, where  $d$  and  $h$  denote the projection size and number of attention heads. Further details are provided in Appendix C. During fine-tuning, all model parameters are updated, including  $\mathbf{W}^\mu$ ,  $\mathbf{b}^\mu$ ,  $\mathbf{W}^\sigma$ ,  $\mathbf{b}^\sigma$ ,  $\mathbf{w}_1^\alpha$ ,  $\mathbf{w}_2^\alpha$ , and  $b^\alpha$ .

To fine-tune with NVIB regularisation, we add Kullback-Leibler (KL) divergence terms to the task loss. As with previous VIB regularisers, information flow is controlled during training by sampling the latent representations. Minimising the KL divergence with the prior tries to maintain this sampling noise and remove information, while the task loss keeps the information needed for the task. The task loss  $\mathcal{L}_T$  is computed with the sampled representations. With NVIB, the KL divergence is decomposed into two loss terms:  $\mathcal{L}_G$  for the Gaussians and  $\mathcal{L}_D$  for the Dirichlet distributions, with hyperparameters  $\lambda_G$  and  $\lambda_D$  controlling their balance. The corresponding equations from Henderson & Fehr (2023) are provided in Appendix A.3. This gives us a total fine-tuning loss of:

$$\mathcal{L} = \mathcal{L}_T + \lambda_D \mathcal{L}_D + \lambda_G \mathcal{L}_G \quad (1)$$

**Novel methods for NVIB fine-tuning.** Firstly, in contrast to Fehr & Henderson (2024), we simplify the denoising function during evaluation to better align with the training function. The equations used in both training and evaluation are shown in Figure 1 (right) and pseudocode in Appendix B. Secondly, while Fehr & Henderson (2024) estimate the prior parameters from training data, in this work we allow the prior mean  $\boldsymbol{\mu}^p$  to be fine-tuned. This allows for flexibility and adaptation to the pretrained model. To maintain the noise in the prior during training, we keep the prior variance  $(\boldsymbol{\sigma}^p)^2 = \mathbf{1}$  and the prior’s pseudo-count  $\alpha_0^p = 1$  fixed. Thirdly, we stabilise fine-tuning by applying proportional clipping to the Dirichlet sampling parameters  $\boldsymbol{\alpha}$ . The magnitude of  $\boldsymbol{\alpha}$  controls the amount of noise when sampling the weights  $\boldsymbol{\pi}$ , with larger values reducing noise. The relative values of  $\boldsymbol{\alpha}$  determine the expected  $\boldsymbol{\pi}$  distribution. Thus, we control the magnitude of  $\boldsymbol{\alpha}$  while preserving its relative values using the clipping functions  $\max(\epsilon, \cdot)$  and  $\min(\omega, \cdot)$  to prevent underflow and overflow, respectively. The parameter  $\epsilon$  is set small enough to prevent values from vanishing, while  $\omega$  is chosen to be sufficiently large to avoid distorting the distribution.

$$\boldsymbol{\alpha} = \max\left(\epsilon, \frac{\boldsymbol{\alpha}}{\sum_i \alpha_i}\right) \times \min\left(\omega, \sum_i \alpha_i\right) \quad (2)$$

### 3 EXPERIMENTS

To evaluate the NVIB regulariser, we design controlled experiments by fine-tuning pretrained models across modalities, including speech, text, graphs, and vision. We compare to models that are first pretrained and then fine-tuned using empirical risk minimization (ERM) with task-specific loss functions. For simplicity and to maintain uniformity across experiments, we define a set of fine-tuned baselines, avoiding modality-specific alternatives. These baselines include models trained without regularisation and models with dropout regularisation. Dropout is a suitable baseline for NVIB regularisation, as it is widely used and effective, seamlessly integrates into pretrained models, and introduces noise into both embeddings and attention mechanisms. To reduce computational costs, we prioritise smaller Transformer models: TinyBERT (Turc et al., 2019) for text and graphs, Wav2Vec2 base and large (Baeovski et al., 2020) for speech, and DeiT-small (Touvron et al., 2021b) for vision. Additional modelling details and hyperparameters for each experiment are provided in Appendix C.

#### 3.1 SPEECH OUT-OF-DISTRIBUTION EVALUATION

Language identification and automated assessment of speech are crucial tasks in the development of audio transmission systems, but are challenging due to many factors related to: the acoustic environment; variation in recording hardware and software; speaker characteristics; and evaluation conditions (Gierlich & Kettler, 2006; Chinen, 2021; Cooper et al., 2022). The prediction of perceived speech quality is formulated as a regression task to estimate the scores of human listeners (ITU-T, 1996), whereas language identification is a classification task given an audio sample. Given the diverse array of factors that can impact speech, generalisation is essential in these tasks.

**Speech quality assessment.** We fine-tune and evaluate on the NISQA (Mittag et al., 2021) dataset, which contains English speech recordings from live calls with network impairments and simulated distortions. We perform OOD testing on the TencentWithReverberation (Tencent) Chinese speech corpus (Yi et al., 2022), which introduces new conditions such as: simulated and real reverberation; and different labelling conditions. Following ITU-T (2020), we evaluate our models using the Pearson’s correlation coefficient (PCC) and root-mean-square error after mapping with a first-order polynomial function (RMSE MAP). Table 1 shows that NVIB regularisation achieves the highest correlation on the in-distribution (ID) data. On the OOD dataset, NVIB regularisation achieves comparable generalisation improvements while exhibiting a lower standard deviation.

Table 1: Speech quality assessment for NISQA (ID) and Tencent (OOD). Average test results (0–1) are reported with standard deviation across 5 seeds.

Model	NISQA (ID)		Tencent (OOD)	
	PCC (↑)	RMSE MAP (↓)	PCC (↑)	RMSE MAP (↓)
W2V2 <sub>Base</sub>	0.89 (0.02)	0.42 (0.03)	0.80 (0.01)	0.54 (0.01)
with Dropout	0.89 (0.01)	0.43 (0.01)	<b>0.83</b> (0.03)	<b>0.51</b> (0.04)
with NVIB	<b>0.90</b> (0.01)	<b>0.41</b> (0.02)	<b>0.83</b> (0.02)	<b>0.51</b> (0.03)

**Speech language identification.** We fine-tune our models on the CommonLanguage (Ravanelli et al., 2021) speech dataset which consists of 22K training audios from 45 languages. We evaluate on two OOD datasets with overlapping languages: FLEURS (Conneau et al., 2023) with 27 languages; and VoxPopuli (Wang et al., 2021) with 11 languages. The FLEURS dataset is read speech, which is closer to CommonLanguage. Whereas, the VoxPopuli dataset is more challenging as it contains spontaneous speech from the European Parliament. Table 2 reports the F1 classification scores, showing that NVIB matches ID performance and outperforms the dropout-regularised baseline on the OOD datasets.

Table 2: Language identification for CommonLanguage (ID), FLEURS (OOD) and VoxPopuli (OOD). Average test F1 scores (0–1) are reported with standard deviation across 5 seeds.

Model	CommonLanguage (ID)	FLEURS (OOD)	VoxPopuli (OOD)
	F1 ( $\uparrow$ )	F1 ( $\uparrow$ )	F1 ( $\uparrow$ )
W2V2 <sub>Large</sub>	<b>0.82</b> (0.01)	0.90 (0.02)	<b>0.86</b> (0.02)
with Dropout	0.81 (0.01)	0.90 (0.01)	0.82 (0.02)
with NVIB	<b>0.82</b> (0.01)	<b>0.91</b> (0.02)	0.85 (0.02)

### 3.2 TEXT OUT-OF-DISTRIBUTION CLASSIFICATION

We consider the CivilComments (CC) (Borkan et al., 2019) task which is part of the WILDS (Koh et al., 2021) curated set of datasets that represent real-life distribution shifts. CC classifies the presence of toxicity in online comments which is an important task of monitoring internet content. The task is a binary classification task of determining if a comment is toxic and contains a subpopulation shift between 8 demographic identities classes. The subpopulation shift means that the training and test domains overlap, but their relative proportions differ. We measure the generalisation by the accuracy of the lowest performing subpopulation *worst-group* (WG).

Table 3 shows the generalisation improvement of this task through regularisation. On average, NVIB regularisation improves OOD generalisation over the unregularised baseline, though it remains less effective than dropout. However, introducing sparsity in the attention keys based on their attention magnitude, as shown in Figure 2, raises the OOD accuracy of the NVIB model and sustains it across a wide range of sparsity levels. Further inspection of the attention patterns in Appendix Figures 4 & 5 shows an interpretable focus on toxic words as spurious keys are dropped and attention weight is put on the prior token.

Table 3: Text classification on CC train (ID) and test (OOD). Average accuracy (%) is reported across 5 seeds with standard deviation and the *best* OOD model.

Model	CC Train (ID)	CC Test (OOD)
	WG ( $\uparrow$ )	WG ( $\uparrow$ )
BERT <sub>Tiny</sub>	78.12 (14.33) 99.00	49.14 (5.56) 61.03
with Dropout	<b>91.05</b> (1.49) 91.16	<b>60.10</b> (3.11) 63.97
with NVIB	80.12 (10.69) 76.30	55.01 (6.15) 61.03

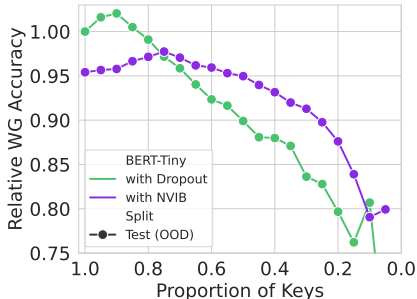


Figure 2: Worst-group (WG) test accuracy as a function of attention key sparsity for the best OOD models, relative to the dropout baseline.

### 3.3 GRAPH LINK PREDICTION

Link prediction is a graph-based problem that involves predicting whether a link exists between two nodes in a graph. This is widely used for recommendation and prediction in social networks, citation links and biological interactions (Kumar et al., 2020; Xia et al., 2021). We build upon the BERT for Link Prediction (BLP) model (Daza et al., 2020) which operates on a set of triples  $(h, r, t)$ , where  $h$  and  $t$  represent the head and tail node, while  $r$  represents the relation between those two nodes. We evaluate on the FB15k-237 dataset (Daza et al., 2020). This dataset follows an inductive setting, where new entities and triples are dynamically incorporated into the graph during evaluation. We evaluate the models by querying them with  $(h, r, ?)$  and  $(?, r, t)$  triples, and assess their performance using two metrics: Mean Reciprocal Rank (MRR), which measures the model’s ability to rank the correct triple, and  $H@k$ , which calculates the proportion of correct triples ranked within the top- $k$  results. Table 4 presents the test set results, which highlights the advantage of the NVIB-regularised model over typical regularisation methods like dropout. This advantage may stem from the presence of new entities in the head or tail positions, which require a higher level of generalisation.

Table 4: Graph link prediction on FB15k-237. Test set ranking metrics (0–1) are reported, based on the best model selected from validation set performance.

Model	FB15k-237			
	MRR ( $\uparrow$ )	H@1 ( $\uparrow$ )	H@3 ( $\uparrow$ )	H@10 ( $\uparrow$ )
BLP-BERT <sub>Tiny</sub>	0.164	0.100	0.175	0.288
with Dropout	0.162	0.097	0.172	0.288
with NVIB	<b>0.167</b>	<b>0.103</b>	<b>0.180</b>	<b>0.294</b>

### 3.4 IMAGE FEW-SHOT CLASSIFICATION

Few-shot classification aims to train models capable of classifying images with limited labelled examples per category. Meta-learning (Vinyals et al., 2016) achieves this by meta-training on several *episodes*, enabling generalisation to new tasks with previously unseen classes. To generalise effectively, the classifier must transfer knowledge from the training distribution to unseen testing distributions while avoiding spurious correlations and shortcuts (Zheng et al., 2024; Zhang et al., 2024). The following experiments are conducted within a meta-learning-based few-shot classification framework Hu et al. (2022), using a pretrained DeiT-small model as the backbone.

**Few-shot in-distribution.** We evaluate the ID performance using the CIFAR-FS (Bertinetto et al., 2019) dataset. Following Hu et al. (2022), we conduct experiments under a 5-way 5-shot setting, where each episode consists of a “support set” with 5 classes and 5 samples per class for training, and a “query set” containing 5 classes with 15 examples per class for testing. Table 5 reports the average classification accuracy and standard deviation over all test episodes for CIFAR-FS in few-shot classification. Compared to the baseline and Dropout, we observe that NVIB regularisation improves accuracy with lower variance across all test episodes.

Table 5: Image classification on CIFAR-FS (ID). Test episodes accuracy (%) with standard deviation.

Model	CIFAR-FS (ID)	
	Acc ( $\uparrow$ )	Std ( $\downarrow$ )
DeiT <sub>Small</sub>	93.57	5.71
with Dropout	93.55	5.61
with NVIB	<b>93.88</b>	<b>5.58</b>

**Few-shot out-of-distribution.** To evaluate the OOD few-shot classification performance, we use the Meta-Dataset (Triantafillou et al., 2019). This benchmark is a diverse set of 10 image datasets, including, ImageNet-1k, MSCOCO (COCO), Traffic Signs (Sign), Describable Textures (DTD), FGVCx Fungi (Fungi), Omniglot, VGG Flower (Flower), CUB-200-2011 (CUB), FGVAircraft (Acraft) and QuickDraw (QDraw). We meta-train the models on ImageNet-1k and then meta-test them on the remaining datasets. The number of ways, shots, and query images for each dataset are sampled as in Hu et al. (2022), with further details provided in Appendix C.4.2. Figure 3 shows that the NVIB-regularised model achieves the highest performance on 6 out of 9 OOD datasets and is rarely outperformed by the dropout-regularised baseline.

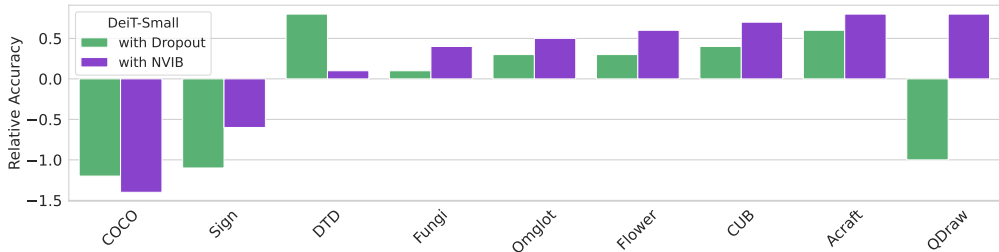


Figure 3: Percentage point improvement in test accuracy relative to the unregularised baseline on the Meta-Dataset benchmark (OOD).

## 4 DISCUSSION

Our results suggest that including NVIB regularisation improves the model’s ability to distinguish signal from noise. This is supported by performance gains observed in tasks such as speech quality prediction (Table 1) and few-shot image classification (Table 5). We attribute this to NVIB’s Bayesian nature, which effectively models statistical uncertainty. During fine-tuning, NVIB introduces noise into the latent representations, which enhances its ability to generalise across noisy feature spaces such as background disturbances and capture variations present in both audio and images. NVIB regularisation shifts the model’s attention from relying on superficial, spurious features to deeper features which generalise better out-of-distribution. This is evident in consistent improvements across tasks that require generalisation to unseen entities, such as graph linking (Table 4) and visual meta-learning (Figure 3). We believe this is due to the additional prior tokens, which disentangle and reweight attention away from spurious tokens (attention maps in Appendix Figures 4 & 5). Additionally, this effect is observed in sustained performance with increased sparsity (Figure 2).

## 5 CONCLUSION

In this work, we contribute to fine-tuning with Nonparametric Variational Information Bottleneck regularisation by demonstrating improved generalisation across multiple modalities and models. We extend NVIB to pretrained models by proposing a novel learnable prior mean embedding per layer for greater adaptability, clipping Dirichlet pseudo-counts for training stability, and simplifying the NVIB denoising attention function at evaluation time.

**Future work.** In future work, we aim to scale our experiments to include models with larger parameter sizes and explore training from scratch. While our current focus prioritized simplicity and uniformity, we are encouraged to evaluate additional baselines and tasks across each modality. Furthermore, we see significant promise in applying NVIB to language modelling, particularly with large language models (LLMs).

## REFERENCES

- Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. Deep variational information bottleneck. In *5th International Conference on Learning Representations, ICLR 2017, Conference Track Proceedings*, Toulon, France, 2017. OpenReview.net. URL <https://openreview.net/forum?id=HyxQzBceg>.
- Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460, 2020.
- Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEit: BERT pre-training of image transformers. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=p-BhZSz59o4>.
- Melika Behjati, Fabio James Fehr, and James Henderson. Learning to abstract with nonparametric variational information bottleneck. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023. URL <https://openreview.net/forum?id=vU0KbvQ91x>.
- Luca Bertinetto, Joao F. Henriques, Philip Torr, and Andrea Vedaldi. Meta-learning with differentiable closed-form solvers. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=HyxnZh0ct7>.
- Prajjwal Bhargava, Aleksandr Drozd, and Anna Rogers. Generalization in NLI: Ways (not) to go beyond simple heuristics. In João Sedoc, Anna Rogers, Anna Rumshisky, and Shabnam Tafreshi (eds.), *Proceedings of the Second Workshop on Insights from Negative Results in NLP*, pp. 125–135, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.insights-1.18. URL <https://aclanthology.org/2021.insights-1.18>.

- Srinadh Bhojanapalli, Ayan Chakrabarti, Andreas Veit, Michal Lukasik, Himanshu Jain, Frederick Liu, Yin-Wen Chang, and Sanjiv Kumar. Leveraging redundancy in attention with reuse transformers. *ArXiv*, abs/2110.06821, 2021. URL <https://api.semanticscholar.org/CorpusID:238743891>.
- Yuchen Bian, Jiaji Huang, Xingyu Cai, Jiahong Yuan, and Kenneth Ward Church. On attention redundancy: A comprehensive study. In *North American Chapter of the Association for Computational Linguistics*, 2021. URL <https://api.semanticscholar.org/CorpusID:235097467>.
- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger (eds.), *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013. URL [https://proceedings.neurips.cc/paper\\_files/paper/2013/file/1cecc7a77928ca8133fa24680a88d2f9-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2013/file/1cecc7a77928ca8133fa24680a88d2f9-Paper.pdf).
- Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Nuanced metrics for measuring unintended bias with real data for text classification. In *Companion Proceedings of The 2019 World Wide Web Conference*, 2019.
- Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *ArXiv*, abs/1904.10509, 2019. URL <https://api.semanticscholar.org/CorpusID:129945531>.
- Michael Chinen. Marginal effects of language and individual raters on speech quality models. *IEEE Access*, 9:127320–127334, 2021. doi: 10.1109/ACCESS.2021.3112165.
- Weiqin Chuah, Ruwan Tennakoon, Reza Hoseinnezhad, Alireza Bab-Hadiashar, and David Suter. Itsa: An information-theoretic approach to automatic shortcut avoidance and domain generalization in stereo matching networks. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 13012–13022, 2022. URL <https://api.semanticscholar.org/CorpusID:245827816>.
- Sanghyuk Chun. Improved probabilistic image-text representations. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=ft1mr3WlGM>.
- Alexis Conneau, Min Ma, Simran Khanuja, Yu Zhang, Vera Axelrod, Siddharth Dalmia, Jason Riesa, Clara Rivera, and Ankur Bapna. Fleurs: Few-shot learning evaluation of universal representations of speech. In *2022 IEEE Spoken Language Technology Workshop (SLT)*, pp. 798–805. IEEE, 2023.
- Erica Cooper, Wen-Chin Huang, Tomoki Toda, and Junichi Yamagishi. Generalization ability of mos prediction networks. In *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 8442–8446, 2022. doi: 10.1109/ICASSP43922.2022.9746395.
- Daniel Daza, Michael Cochez, and Paul T. Groth. Inductive entity representations from text via link prediction. *Proceedings of the Web Conference 2021*, 2020. URL <https://api.semanticscholar.org/CorpusID:222177425>.
- Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, Rodolphe Jenatton, Lucas Beyer, Michael Tschannen, Anurag Arnab, Xiao Wang, Carlos Riquelme Ruiz, Matthias Minderer, Joan Puigcerver, Utku Evci, Manoj Kumar, Sjoerd Van Steenkiste, Gamaleldin Fathy Elsayed, Aravindh Mahendran, Fisher Yu, Avital Oliver, Fantine Huot, Jasmijn Bastings, Mark Collier, Alexey A. Gritsenko, Vighnesh Birodkar, Cristina Nader Vasconcelos, Yi Tay, Thomas Mensink, Alexander Kolesnikov, Filip Pavetic, Dustin Tran, Thomas Kipf, Mario Lucic, Xiaohua Zhai, Daniel Keysers, Jeremiah J. Harmsen, and Neil Houlsby. Scaling vision transformers to 22 billion parameters. In Andreas Krause, Emma



- Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 7480–7512. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/dehghani23a.html>.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- William Fedus, Barret Zoph, and Noam M. Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *J. Mach. Learn. Res.*, 23:120:1–120:39, 2021. URL <https://api.semanticscholar.org/CorpusID:231573431>.
- Fabio James Fehr and James Henderson. Nonparametric variational regularisation of pre-trained transformers. In *First Conference on Language Modeling*, 2024. URL <https://openreview.net/forum?id=Zu8OWNUC0u>.
- R. Geirhos, J.-H. Jacobsen, C. Michaelis, R. Zemel, W. Brendel, M. Bethge, and F. A. Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020. doi: 10.1038/s42256-020-00257-z.
- H.W. Gierlich and F. Kettler. Advanced speech quality testing of modern telecommunication equipment: An overview. *Signal Processing*, 86(6):1327–1340, 2006. ISSN 0165-1684. doi: <https://doi.org/10.1016/j.sigpro.2005.06.024>. URL <https://www.sciencedirect.com/science/article/pii/S0165168405003312>. Applied Speech and Audio Processing.
- Zongyan Han, Zhenyong Fu, and Jian Yang. Learning the redundancy-free features for generalized zero-shot object recognition. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12862–12871, 2020. URL <https://api.semanticscholar.org/CorpusID:219633468>.
- James Henderson. The unstoppable rise of computational linguistics in deep learning. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 6294–6306, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.561. URL <https://aclanthology.org/2020.acl-main.561>.
- James Henderson and Fabio James Fehr. A VAE for transformers with nonparametric variational information bottleneck. In *The Eleventh International Conference on Learning Representations*, 2023. URL [https://openreview.net/forum?id=6QkjC\\_cs03X](https://openreview.net/forum?id=6QkjC_cs03X).
- Shell Xu Hu, Da Li, Jan Stühmer, Minyoung Kim, and Timothy M. Hospedales. Pushing the limits of simple pipelines for few-shot learning: External data and fine-tuning make a difference. In *CVPR*, 2022.
- ITU-T. Methods for subjective determination of transmission quality. Recommendation P.800, International Telecommunication Union, Geneva, Switzerland, August 1996.
- ITU-T. Methods, metrics and procedures for statistical evaluation, qualification and comparison of objective quality prediction models. Recommendation P.1401, International Telecommunication Union, Geneva, Switzerland, January 2020.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014. URL <https://api.semanticscholar.org/CorpusID:6628106>.

- Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, Tony Lee, Etienne David, Ian Stavness, Wei Guo, Berton A. Earnshaw, Imran S. Haque, Sara Beery, Jure Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang. Wilds: A benchmark of in-the-wild distribution shifts, 2021.
- Ajay Kumar, Shashank Sheshar Singh, Kuldeep Singh, and Bhaskar Biswas. Link prediction techniques, applications, and performance: A survey. *Physica A: Statistical Mechanics and its Applications*, 553:124289, 2020. ISSN 0378-4371. doi: <https://doi.org/10.1016/j.physa.2020.124289>. URL <https://www.sciencedirect.com/science/article/pii/S0378437120300856>.
- Ananya Kumar, Aditi Raghunathan, Robbie Matthew Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=UYneFzXSJWh>.
- Bo Li, Yifei Shen, Yezhen Wang, Wenzhen Zhu, Colorado Reed, Jun Zhang, Dongsheng Li, Kurt Keutzer, and Han Zhao. Invariant information bottleneck for domain generalization. In *AAAI Conference on Artificial Intelligence*, 2021a. URL <https://api.semanticscholar.org/CorpusID:235417355>.
- Pengyong Li, Jun Wang, Yixuan Qiao, Hao Chen, Yihuan Yu, Xiaojun Yao, Peng Gao, Guowang Xie, and Sen Song. An effective self-supervised framework for learning expressive molecular global representations to drug discovery. *Briefings in bioinformatics*, 2021b. URL <https://api.semanticscholar.org/CorpusID:233719686>.
- Jiachen Lian, Chunlei Zhang, and Dong Yu. Robust disentangled variational speech representation learning for zero-shot voice conversion. *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6572–6576, 2022. URL <https://api.semanticscholar.org/CorpusID:247839160>.
- Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rkgz2aEKDr>.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 10012–10022, October 2021.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- Rabeeh Karimi mahabadi, Yonatan Belinkov, and James Henderson. Variational information bottleneck for effective low-resource fine-tuning. In *International Conference on Learning Representations*, 2021. URL [https://openreview.net/forum?id=kvhzKz-\\_DMF](https://openreview.net/forum?id=kvhzKz-_DMF).
- Arya D. McCarthy, Xian Li, Jiatao Gu, and Ning Dong. Addressing posterior collapse with mutual information for improved variational neural machine translation. In *Annual Meeting of the Association for Computational Linguistics*, 2020. URL <https://api.semanticscholar.org/CorpusID:220046608>.
- Gabriel Mittag, Babak Naderi, Assmaa Chehadi, and Sebastian Möller. Nisqa: A deep cnn-self-attention model for multidimensional speech quality prediction with crowdsourced datasets. *Interspeech 2021*, Aug 2021. doi: 10.21437/interspeech.2021-299. URL <http://dx.doi.org/10.21437/Interspeech.2021-299>.
- Alexandru Nelus and Rainer Martin. Privacy-preserving audio classification using variational information feature extraction. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:2864–2877, 2021. URL <https://api.semanticscholar.org/CorpusID:237518792>.

- Bhargavi Paranjape, Mandar Joshi, John Thickstun, Hannaneh Hajishirzi, and Luke Zettlemoyer. An information bottleneck approach for controlling conciseness in rationale extraction. In *Conference on Empirical Methods in Natural Language Processing*, 2020. URL <https://api.semanticscholar.org/CorpusID:218487373>.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, pp. 28492–28518. PMLR, 2023.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL <http://jmlr.org/papers/v21/20-074.html>.
- Mirco Ravanelli, Titouan Parcollet, Peter Plantinga, Aku Rouhe, Samuele Cornell, Loren Lugosch, Cem Subakan, Nauman Dawalatabad, Abdelwahab Heba, Jianyuan Zhong, et al. Speechbrain: A general-purpose speech toolkit. *arXiv preprint arXiv:2106.04624*, 2021.
- Yu Rong, Yatao Bian, Tingyang Xu, Weiyang Xie, Ying WEI, Wenbing Huang, and Junzhou Huang. Self-supervised graph transformer on large-scale molecular data. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 12559–12571. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/94aef38441efa3380a3bed3faf1f9d5d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/94aef38441efa3380a3bed3faf1f9d5d-Paper.pdf).
- Sebastian Ruder, Matthew E. Peters, Swabha Swayamdipta, and Thomas Wolf. Transfer learning in natural language processing. In *North American Chapter of the Association for Computational Linguistics*, 2019. URL <https://api.semanticscholar.org/CorpusID:186206211>.
- Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems*, 2017.
- Qingyun Sun, Jianxin Li, Hao Peng, Jia Wu, Xingcheng Fu, Cheng Ji, and Philip S. Yu. Graph structure learning with variational information bottleneck. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(4):4165–4174, 2022. ISSN 2374-3468. doi: 10.1609/aaai.v36i4.20335. 36th AAAI Conference on Artificial Intelligence, AAAI 2022 ; Conference date: 22-02-2022 Through 01-03-2022.
- Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. *2015 IEEE Information Theory Workshop (ITW)*, pp. 1–5, 2015. URL <https://api.semanticscholar.org/CorpusID:5541663>.
- Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *ArXiv*, physics/0004057, 2000. URL <https://api.semanticscholar.org/CorpusID:8936496>.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. Training data-efficient image transformers and distillation through attention. In *International Conference on Machine Learning*, volume 139, pp. 10347–10357, July 2021a.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. Training data-efficient image transformers amp; distillation through attention. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 10347–10357. PMLR, 18–24 Jul 2021b. URL <https://proceedings.mlr.press/v139/touvron21a.html>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *ArXiv*, abs/2302.13971, 2023. URL <https://api.semanticscholar.org/CorpusID:257219404>.

- Eleni Triantafyllou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Utku Evci, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, et al. Meta-dataset: A dataset of datasets for learning to learn from few examples. In *International Conference on Learning Representations*, 2019.
- Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Well-read students learn better: On the importance of pre-training compact models. *arXiv: Computation and Language*, 2019. URL <https://api.semanticscholar.org/CorpusID:202889175>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. *Advances in neural information processing systems*, 29, 2016.
- Changhan Wang, Morgane Riviere, Ann Lee, Anne Wu, Chaitanya Talnikar, Daniel Haziza, Mary Williamson, Juan Pino, and Emmanuel Dupoux. VoxPopuli: A large-scale multilingual speech corpus for representation learning, semi-supervised learning and interpretation. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 993–1003, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.80. URL <https://aclanthology.org/2021.acl-long.80>.
- Jiawei Wu, Xiaoya Li, Xiang Ao, Yuxian Meng, Fei Wu, and Jiwei Li. Improving robustness and generality of nlp models using disentangled representations. *ArXiv*, abs/2009.09587, 2020a. URL <https://api.semanticscholar.org/CorpusID:221819589>.
- Tailin Wu, Hongyu Ren, Pan Li, and Jure Leskovec. Graph information bottleneck. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 20437–20448. Curran Associates, Inc., 2020b. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/ebc2aa04e75e3caabda543a1317160c0-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/ebc2aa04e75e3caabda543a1317160c0-Paper.pdf).
- Feng Xia, Ke Sun, Shuo Yu, Abdul Aziz, Liangtian Wan, Shirui Pan, and Huan Liu. Graph learning: A survey. *IEEE Transactions on Artificial Intelligence*, 2:109–127, 2021. URL <https://api.semanticscholar.org/CorpusID:233481068>.
- Gaoxiong Yi, Wei Xiao, Yiming Xiao, Babak Naderi, Sebastian Möller, Wafaa Wardah, Gabriel Mittag, Ross Culter, Zhuohuang Zhang, Donald S. Williamson, Fei Chen, Fuzheng Yang, and Shidong Shang. ConferencingSpeech 2022 Challenge: Non-intrusive Objective Speech Quality Assessment (NISQA) Challenge for Online Conferencing Applications. In *Proc. Interspeech 2022*, pp. 3308–3312, 2022. doi: 10.21437/Interspeech.2022-10597.
- Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12104–12113, June 2022.
- Cenyuan Zhang, Xiang Zhou, Yixin Wan, Xiaoqing Zheng, Kai-Wei Chang, and Cho-Jui Hsieh. Improving the adversarial robustness of NLP models by information bottleneck. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 3588–3598, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.284. URL <https://aclanthology.org/2022.findings-acl.284/>.
- Min Zhang, Haoxuan Li, Fei Wu, and Kun Kuang. Metacoco: A new few-shot classification benchmark with spurious correlation. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=DiWRG9JTWZ>.

Guangtao Zheng, Wenqian Ye, and Aidong Zhang. Benchmarking spurious bias in few-shot image classifiers. In *European Conference on Computer Vision*, 2024. URL <https://api.semanticscholar.org/CorpusID:272398061>.

## AUTHOR CONTRIBUTIONS & ACKNOWLEDGMENTS

Fabio Fehr lead the project and integrated NVIB into all pretrained models while conducting experiments on OOD text classification. He was supported by the Swiss National Centre of Competence in Research (NCCR) under the project Evolving Language, grant number “51NF40\_180888”. Alina Elena Baia worked on initial experimentation of NVIB for vision including robustness and privacy classification and contributed to the writing of the paper. Xiaoguang Chang worked on few-shot image classification experiments and analysis. Andrei C. Coman worked on graph link prediction and improvements to the architecture. He was supported by the Swiss National Science Foundation (SNSF) under the project “Deep Learning Models for Continual Extraction of Knowledge from Text”, grant number “200021E\_189458”. Karl El Hajal performed experiments on speech quality. He was supported by the Swiss National Science Foundation grant agreement no. 219726 on “Pathological Speech Synthesis (PaSS)”. Dina El Zein contributed to initial, exploratory experiments on NVIB for text-related tasks. Shashi Kumar performed experiments on speech language identification. Juan Zuluaga-Gomez worked on initial experimentation of NVIB for speech related tasks. Andrea Cavallaro provided guidance, supervision and contributed to the writing of the paper. Damien Teney provided guidance and helped editing the paper. James Henderson provided the primary guidance and supervision while contributing to writing the theoretical sections of the paper.

## A INTRODUCTION TO NVIB

Henderson & Fehr (2023) define Nonparametric Variational Information Bottleneck (NVIB) by generalising the standard attention layer to a Bayesian model where embeddings are distributions over the latent space. A key insight of this approach is that the latent space of attention-based representations can be viewed as non-parametric mixture distributions. In this interpretation, the vectors accessed via attention define a mixture of impulse distributions. Since a Transformer embedding is a set of vectors that dynamically scale with the complexity of the input, the corresponding latent space of these mixture distributions is inherently nonparametric in nature. In this formulation, the attention function is interpreted as Bayesian “query denoising” using the latent distribution as the prior. The authors define *denoising attention* as a generalisation of the attention function to query denoising.

### A.1 DENOISING ATTENTION

Denoising attention is a generalisation of attention which interprets the latent space of Transformers as a non-parametric mixture distribution. Henderson & Fehr (2023) provide a constructive proof of exact equivalence to the standard attention function. When standard attention accesses the latent space of Transformers, which is a set of embedding vectors  $\mathbf{Z} \in \mathbb{R}^{n \times d}$  via weight matrices  $\mathbf{W}^K, \mathbf{W}^V \in \mathbb{R}^{d \times d}$  to keys and values, respectively, and projects the accessing input vector  $\mathbf{u}' \in \mathbb{R}^{1 \times d}$  via the weight matrix  $\mathbf{W}^Q \in \mathbb{R}^{d \times d}$  to a query. By letting  $\mathbf{u} = (\mathbf{u}' \mathbf{W}^Q (\mathbf{W}^K)^\top) \in \mathbb{R}^{1 \times d}$ , the standard scaled dot product attention function can be rewritten as  $(\text{Attn}(\mathbf{u}, \mathbf{Z}) \mathbf{W}^V)$ , with  $\text{Attn}(\mathbf{u}, \mathbf{Z})$  defined in terms of a sum over the vectors  $\mathbf{z}_i$  in  $\mathbf{Z}$ , or equivalently defined in terms of an integral over a distribution which is only non-zero at the  $\mathbf{z}_i$ :

$$\text{Attn}(\mathbf{u}, \mathbf{Z}) = \text{softmax} \left( \frac{1}{\sqrt{d}} \mathbf{u} \mathbf{Z}^\top \right) \mathbf{Z} = \text{DAttn}(\mathbf{u}; F_{\mathbf{Z}}) \quad (3)$$

$$\text{DAttn}(\mathbf{u}; F) = \int_{\mathbf{v}} \frac{f(\mathbf{v}) g(\mathbf{u}; \mathbf{v}, \sqrt{d} \mathbf{I})}{\int_{\mathbf{v}} f(\mathbf{v}) g(\mathbf{u}; \mathbf{v}, \sqrt{d} \mathbf{I}) d\mathbf{v}} \mathbf{v} d\mathbf{v} \quad (4)$$

$$F_{\mathbf{Z}} = \sum_{i=1}^n \frac{\exp(\frac{1}{2\sqrt{d}} \|\mathbf{z}_i\|^2)}{\sum_{i=1}^n \exp(\frac{1}{2\sqrt{d}} \|\mathbf{z}_i\|^2)} \delta_{\mathbf{z}_i} \quad (5)$$

where  $\delta_{\mathbf{z}_i}$  is an impulse distribution at  $\mathbf{z}_i$ ,  $f(\cdot)$  is the probability density function for distribution  $F$ , and  $g(\mathbf{u}; \mathbf{v}, \sqrt{d} \mathbf{I})$  is the multivariate Gaussian function with diagonal variance of  $\sqrt{d}$ . This alternative definition  $\text{DAttn}(\mathbf{u}; F_{\mathbf{Z}})$  is *denoising attention*. It subsumes standard attention in that any attention-based embedding  $\mathbf{Z}$  has an equivalent mixture of impulse distributions, namely  $F_{\mathbf{Z}}$ , where denoising attention  $\text{DAttn}(\mathbf{u}; F_{\mathbf{Z}})$  gives us exactly the same result as attention  $\text{Attn}(\mathbf{u}, \mathbf{Z})$ , for all queries  $\mathbf{u}$ . This is an elegant result, which in

practice allows us to define a nonparametric distribution over the latent embeddings of Transformers. Appendix B covers the exact equations for denoising attention and pseudocode.

## A.2 DISTRIBUTIONS OVER MIXTURE DISTRIBUTIONS

Given this generalisation of attention-based representations to nonparametric mixture distributions, Bayesian nonparametrics can be used to define distributions over the latent space. Henderson & Fehr (2023) propose to use Dirichlet Processes (DPs) to define distributions over mixture distributions, so an NVIB layer first embeds its input vectors into a DP representation by mapping them to the parameters  $(\boldsymbol{\mu}^q, \boldsymbol{\sigma}^q, \boldsymbol{\alpha}^q)$  of a DP. A DP is defined by a base distribution  $G_0^q$  for generating the vectors for the component impulse distributions, and a pseudo-count  $\alpha_0^q$  for generating their mixture weights.

$$\alpha_0^q = \sum_i \alpha_i^q ; \quad G_0^q = \sum_i \frac{\alpha_i^q}{\alpha_0^q} \mathcal{N}(\boldsymbol{\mu}_i^q, \mathbf{I}(\boldsymbol{\sigma}_i^q)^2) \quad (6)$$

Following this definition,  $G_0^q$  is itself a mixture distribution, consisting of one Gaussian component from the prior plus one Gaussian component for each vector input to the NVIB layer. These DPs represent the posterior  $q(F|x)$ . The prior  $p(F)$  is a DP specified by the parameters  $(\boldsymbol{\mu}^p, \boldsymbol{\sigma}^p, \alpha^p)$  of its pseudo-count  $\alpha^p$  and its unimodal base distribution  $G_0^p = \mathcal{N}(\boldsymbol{\mu}^p, \boldsymbol{\sigma}^p)$ . In this work, we allow the prior  $\boldsymbol{\mu}^p$  to be learned, which allows the prior to be centred in the latent embedding space. However, to maintain noise during regularisation, we set the prior variance  $(\boldsymbol{\sigma}^p)^2 = \mathbf{1}$  and the prior’s pseudo-count  $\alpha_0^p = 1$ .

## A.3 NVIB REGULARISATION

During training, NVIB regularises the information passing through the NVIB layer by sampling latent representations from its DP embedding. This process introduces noise and removes redundant information, enhancing model generalisation. The level of noise is learned by the DP parameters  $(\boldsymbol{\mu}^q, \boldsymbol{\sigma}^q, \boldsymbol{\alpha}^q)$  within the NVIB layer. To maintain noise during training, a Kullback-Leibler (KL) divergence loss term is included between the embedding distribution and the DP prior. Since the prior DP is input independent, the KL term enforces an information bottleneck by minimising the information retained in the DP embedding. During evaluation, the NVIB layer uses the mean latent representation, which is the base distribution  $G_0^q$  of the DP embedding.

The evidence lower bound (ELBO) is a widely used objective in variational Bayesian methods, serving as a tractable approximation to the log-likelihood of the observation  $x$ , where  $x$  represents the input. The ELBO is formulated as follows:

$$\log(p(x)) \geq \mathbb{E}_{q(F|x)} \log(p(x|F)) - \mathbb{KL}(q(F|x)||p(F)) \quad (7)$$

$$\mathcal{L}_R = -\mathbb{E}_{q(F|x)} \log(p(x|F)) \quad (8)$$

This decomposition consists of two key terms: the reconstruction loss  $\mathcal{L}_R$ , computed using samples  $F$  drawn from the approximate posterior  $q(F|x)$ , and the KL divergence between this posterior and the prior  $p(F)$ . In this work, we replace the reconstruction loss with a task specific loss  $\mathcal{L}_T$ . Henderson & Fehr (2023) further divided the KL term into  $\mathcal{L}_G$ , corresponding to Gaussian distributions, and  $\mathcal{L}_D$ , corresponding to Dirichlet distributions. This gives us the following loss terms for the KL divergence, where  $\Gamma$  is the gamma function and  $\psi$  is the digamma function:

$$\mathcal{L}_D + \mathcal{L}_G \approx \mathbb{D}_{\mathbb{KL}}(q(F|x) || p(F)) \quad (9)$$

$$\mathcal{L}_D = \log\Gamma(\alpha_0^q) - \log\Gamma(\alpha_0^p) + (\alpha_0^q - \alpha_0^p) \left( -\psi(\alpha_0^q) + \psi\left(\frac{\alpha_0^q}{\kappa_0}\right) \right) + \kappa_0 \left( \log\Gamma\left(\frac{\alpha_0^p}{\kappa_0}\right) - \log\Gamma\left(\frac{\alpha_0^q}{\kappa_0}\right) \right) \quad (10)$$

$$\mathcal{L}_G = \frac{1}{2} \kappa_0 \sum_{i=1}^{n+1} \frac{\alpha_i^q}{\alpha_0^q} \sum_{h=1}^d \left( \frac{(\mu_{ih}^q - \mu_h^p)^2}{(\sigma_h^p)^2} + \frac{(\sigma_{ih}^q)^2}{(\sigma_h^p)^2} - 1 - \log \frac{(\sigma_{ih}^q)^2}{(\sigma_h^p)^2} \right) \quad (11)$$

Since we only draw a single sample per component, thus  $\kappa_0 = n + 1$ . However, in practice we scale both  $\mathcal{L}_G$  and  $\mathcal{L}_D$  by the number of components  $(n + 1)$  such that the loss is invariant to sequence length. We introduce two hyperparameters to control the relative weight of the above three parts of the loss, which defines our VIB loss  $\mathcal{L}$ .

$$\mathcal{L} = \mathcal{L}_T + \lambda_D \mathcal{L}_D + \lambda_G \mathcal{L}_G \quad (12)$$

#### A.4 INCLUDING NVIB INTO PRETRAINED MODELS

Fehr & Henderson (2024) define an identity initialisation for NVIB such that the latent embeddings have negligible uncertainty and denoising attention is effectively equivalent to standard attention. This allows pretrained attention-based models to be reinterpreted as Nonparametric Variational models. By only changing the initialisation, away from the identity and towards an empirically estimated prior, an effective post-training regularisation is added. The authors found that this information-theoretic regularisation lead to improvements in OOD text generalisation in summarisation and translation without fine-tuning.

## B SIMPLIFYING DENOISING ATTENTION

In this section, we provide the implementation details for *denoising multihead attention*. We define the set of Transformer latent embedding vectors as  $\mathbf{Z} \in \mathbb{R}^{n \times d}$  and set of pre-projected queries as  $\mathbf{U}' \in \mathbb{R}^{m \times d}$ . We assume the latent projection matrices are square such that  $\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V \in \mathbb{R}^{d \times d}$  and biases  $\mathbf{b}^Q, \mathbf{b}^K, \mathbf{b}^V \in \mathbb{R}^d$  are used to linearly project to the queries, keys and values, respectively. We define the standard attention weights before the softmax as follows:

$$\mathbf{A} = \frac{1}{\sqrt{d}} \underbrace{(\mathbf{U}'\mathbf{W}^Q + \mathbf{b}^Q)}_{\mathbf{Q}} \underbrace{(\mathbf{Z}\mathbf{W}^K + \mathbf{b}^K)^\top}_{\mathbf{K}^\top} \in \mathbb{R}^{m \times n} \quad (13)$$

Typically, for multihead attention the projected query  $\mathbf{Q}$  and keys  $\mathbf{K}$  are split into heads. In this definition, we split the linear projections by a divisible number of heads  $h$  such that  $\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V \in \mathbb{R}^{h \times d \times \frac{d}{h}}$  and biases  $\mathbf{b}^Q, \mathbf{b}^K, \mathbf{b}^V \in \mathbb{R}^{h \times \frac{d}{h}}$ , so that  $\mathbf{Q} \in \mathbb{R}^{h \times m \times \frac{d}{h}}$  and  $\mathbf{K} \in \mathbb{R}^{h \times n \times \frac{d}{h}}$ . We can then specify multihead attention by defining a matrix of attention scores  $\mathbf{A} \in \mathbb{R}^{h \times m \times n}$ , for each head  $i$ :

$$\mathbf{A}_i = \frac{1}{\sqrt{d/h}} ((\mathbf{Q}_i(\mathbf{W}_i^K)^\top \mathbf{Z}^\top + \mathbf{Q}_i(\mathbf{b}_i^K)^\top) \quad (14)$$

where the bias term  $\mathbf{Q}_i(\mathbf{b}_i^K)^\top \in \mathbb{R}^m$  is added across all  $n$  keys, and thus is normalised out in the softmax below. The scaling term also considers the heads and is division by  $\sqrt{d/h}$ . For denoising attention, each head's query is projected into the space of the original set of vectors  $\mathbf{Z}$ , namely  $\mathbf{U}_i = \mathbf{Q}_i(\mathbf{W}_i^K)^\top$ , and so is still in  $\mathbb{R}^{m \times d}$ . Thus, each head can be viewed as doing denoising attention in the same way as single-head attention, with the only difference being that the variance of the theoretical query noise is now  $\sqrt{d/h}\mathbf{I}$ .

**Training.** Given these sampled weights and vectors, the training-time denoising attention function is the same as the standard attention function with two changes: (1) the keys come from the sampled vectors  $\mathbf{Z} \in \mathbb{R}^{(n+1) \times d}$ , which include a vector sampled from the prior component; and (2) each key has an attention bias  $\mathbf{b} \in \mathbb{R}^{(n+1)}$  which is determined by its weight  $\boldsymbol{\pi} \in \mathbb{R}^{(n+1)}$ . Summing over heads  $i$ , the training-time denoising attention function is:

$$\text{DAttn}(\cdot) = \sum_i \text{Softmax}(\underbrace{\mathbf{A}_i + \log(\boldsymbol{\pi})}_{\mathbf{b}} - \underbrace{\frac{1}{2\sqrt{d/h}}\|\mathbf{Z}\|^2}_{\mathbf{v}_i}) (\mathbf{Z}\mathbf{W}_i^V + \mathbf{b}_i^V) \quad (15)$$

The biases  $\mathbf{b}$  are defined by adding the log of the sampled weights  $\log(\boldsymbol{\pi}) \in \mathbb{R}^{(n+1)}$  from the NVIB layer and subtracting the scaled squared-L2-norms of the sampled vectors  $\frac{1}{2\sqrt{d/h}}\|\mathbf{Z}\|^2 \in \mathbb{R}^{(n+1)}$ . For multihead attention we only need to reuse the same biases  $\mathbf{b}$  for each head, just like we reuse the same vectors  $\mathbf{Z}$  for each head.

**Evaluation.** During the evaluation, as for training, the NVIB layer outputs the isotropic Gaussian parameters  $\boldsymbol{\mu} \in \mathbb{R}^{(n+1) \times d}$ ,  $\boldsymbol{\sigma} \in \mathbb{R}^{(n+1) \times d}$  and Dirichlet parameters  $\boldsymbol{\alpha} \in \mathbb{R}^{(n+1)}$ . For evaluation the base distribution is used. The parameters are taken directly without sampling such that we use the expectation of the distribution. We can write the denoising attention scores  $\mathbf{A} \in \mathbb{R}^{h \times m \times (n+1)}$ , for each head  $i$ , as follows:

$$\mathbf{A}_i = \mathbf{Q}_i(\mathbf{W}_i^K)^\top \left(\frac{\boldsymbol{\mu}}{\sqrt{d/h}}\right)^\top + \frac{1}{\sqrt{d/h}} \mathbf{Q}_i(\mathbf{b}_i^K)^\top \quad (16)$$

where the bias term  $\mathbf{Q}_i(\mathbf{b}_i^K)^\top \in \mathbb{R}^m$  is added across all  $n$  keys, and thus is normalised out in the softmax below. For this attention score matrix  $\mathbf{A}$ , multihead evaluation denoising attention adds the same key biases  $\mathbf{c} \in \mathbb{R}^{h \times (n+1)}$  across all  $m$  queries and  $h$  heads. For ease of notation we define  $\alpha_0 = \sum_{j=1}^d \alpha_j$ .

$$\text{DAttn}(\cdot) = \sum_i \text{Softmax}(\underbrace{\mathbf{A}_i + \log\left(\frac{\boldsymbol{\alpha}}{\alpha_0}\right) - \frac{1}{2}\left\|\frac{\boldsymbol{\mu}}{\sqrt{d}}\right\|^2}_{\mathbf{b}}) \underbrace{(\boldsymbol{\mu}\mathbf{W}_i^V + \mathbf{b}_i^V)}_{\mathbf{v}_i} \quad (17)$$

This simplifies previous implementations of Henderson & Fehr (2023) and Fehr & Henderson (2024) by removing the additional variance term in the bias  $b$  and the interpolation between the query and value vectors. This makes the training and test time denoising attention functions more similar and reduces computation requirements.

Pseudocode: Attention and Denoising Attention during training (single-head). Left: Standard Attention. Right: Denoising Attention.

```
class Attention():
    def __init__(self, d):
        # Projections to Q, K, V [d,d]
        self.q = linear(d, d)
        self.k = linear(d, d)
        self.v = linear(d, d)

    def forward(self, u, z):
        # queries u: [B, M, d]
        # keys / values z: [B, N, d]
        d = keys.shape(2)

        # Project to Q, K, V
        q = self.q(u)
        k = self.k(z) / sqrt(d)
        v = self.v(z)

        # Attention scores [B, M, N]
        attn = q @ k.transpose()

        # Attention probabilities [B, M, N]
        attn = softmax(attn)

        # Value projection [B, M, d]
        out = attn @ v

    return out
```

```
class DenoisingAttention():
    def __init__(self, d):
        # Projections to Q, K, V [d,d]
        self.q = linear(d, d)
        self.k = linear(d, d)
        self.v = linear(d, d)

    def forward(self, u, z, pi):
        # queries u: [B, M, d]
        # keys / values z: [B, N+1, d]
        d = keys.shape(2)

        # Project to Q, K, V
        q = self.q(u)
        k = self.k(z) / sqrt(d)
        v = self.v(z)

        # NVIB bias [B, 1, N+1]
        b = log(pi)
        - 1/(2*sqrt(d))*l2norm(z)**2

        # Attention scores [B, M, N+1]
        attn = q @ k.transpose() + b

        # Attention probabilities [B, M, N+1]
        attn = softmax(attn)

        # Value projection [B, M, d]
        out = attn @ v

    return out
```

Pseudocode: Denoising Attention during evaluation (single-head). Left: Previous implementation including extra bias term and query value interpolation. Right: Current simplified implementation.

```
class DenoisingAttention():
    def __init__(self, d):
        # Projections to Q, K, V [d,d]
        self.q = linear(d, d)
        self.k = linear(d, d)
        self.v = linear(d, d)

    def forward(self, u, mu, sigma2, alpha):
        # queries u: [B, M, d]
        # keys / values mu: [B, N+1, d]
        d = keys.shape(2)

        # Project to Q, K, V
        q = self.q(u)
        k = self.k(mu / (sqrt(d)+sigma2))
        # v is used in interpolation

        # NVIB bias [B, 1, N+1]
        b = log(alpha / sum(alpha))
        - 1/(2*(sqrt(d)+sigma2))*l2norm(mu)**2
        - sum(log(sqrt(sqrt(d)+sigma2)))

        # Attention scores [B, M, N+1]
        attn = q @ k.transpose() + b

        # Attention probabilities [B, M, N+1]
        attn = softmax(attn)

        # Query projection to key-space [B, M, d]
        u_k = self.k(q)

        # Value interpolation [B, M, d]
        out = (attn @
        (sigma2/(sqrt(d)+sigma2))*u_k
        + attn @
        ((sqrt(d)/(sqrt(d)+sigma2))*mu
        out = self.v(out)

    return out
```

```
class DenoisingAttention():
    def __init__(self, d):
        # Projections to Q, K, V [d,d]
        self.q = linear(d, d)
        self.k = linear(d, d)
        self.v = linear(d, d)

    def forward(self, u, mu, alpha):
        # queries u: [B, M, d]
        # keys/values mu: [B, N+1, d]
        d = keys.shape(2)

        # Project to Q, K, V
        q = self.q(u)
        k = self.k(mu) / sqrt(d)
        v = self.v(mu)

        # NVIB bias [B, 1, N+1]
        b = log(alpha/sum(alpha))
        - 1/(2*sqrt(d))*l2norm(mu)**2

        # Attention scores [B, M, N+1]
        attn = q @ k.transpose() + b

        # Attention probabilities [B, M, N+1]
        attn = softmax(attn)

        # Value projection [B, M, d]
        out = attn @ v

    return out
```



## C MODELLING & HYPERPARAMETERS

We outline the general modelling choices applied across all experiments, followed by experiment-specific configurations in the sections below. To manage computational costs, we prioritise smaller models, as all experiments were conducted on a consumer-grade GPU (NVIDIA RTX3090 24GB).

**Baselines.** For uniformity across modalities and models, we use two baselines to compare our regularisation method. The first is an unregularised model fine-tuned without dropout in embeddings and attention. The second, with dropout, follows typical regularisation in pretrained Transformers, applying a 0.1 dropout rate during fine-tuning. Dropout is an appropriate baseline for NVIB regularisation, as it is widely used and effective, seamlessly integrates into pretrained models, and introduces noise into both embeddings and attention mechanisms.

**Initialisation of NVIB layers.** When including the NVIB layers into a pretrained Transformer, we ensure an equivalence in the initialisation, as described by Fehr & Henderson (2024). Specifically, this requires the attention weight to completely ignore the prior component embedding for each layer of the model that includes NVIB. While Fehr & Henderson (2024) empirically initialise the prior component from training data, we simplify the process by initialising our prior mean  $\boldsymbol{\mu}^p = \mathbf{0}$ , variance  $(\boldsymbol{\sigma}^p)^2 = \mathbf{1}$  (standard normal Gaussian), and prior Dirichlet pseudo-count  $\alpha_0^p = 1$ . We establish a set of equivalence tests and find that, in general, the lower layers of the model are easier to preserve in equivalence. The higher layers tend to be more sensitive, and stacking multiple NVIB layers introduces accumulating variance that can break the equivalence. For the later layers where equivalence is not achieved, we exclude NVIB from those layers, with a manual process for each pretrained model. The initialisation Gaussian variance  $\tau_\sigma$  and Dirichlet  $\tau_\alpha$  parameters influence equivalence during both evaluation and training.

The initialisation Gaussian variance parameter  $\tau_\sigma$  is a bias term for the initial amount of noise during fine-tuning. Since it is not required for equivalence, as we do not sample from the embedding distribution during evaluation, we can start with a non-zero amount of noise and initialise this parameter for fine-tuning.

$$\boldsymbol{\sigma}^2 = \sigma^2(\mathbf{x}) = \exp(\mathbf{x}\mathbf{W}^\sigma + \mathbf{b}^\sigma); \quad \mathbf{W}^\sigma = \mathbf{0}; \quad \mathbf{b}^\sigma = \log(\tau_\sigma^2) \quad (18)$$

The initialisation parameter  $\tau_\alpha$  serves as a bias term for the Dirichlet pseudo-counts projection. It reweights the  $\alpha$  parameters, ensuring that the embedding vectors from the input are larger than the prior in the attention calculation. The  $\tau_\alpha$  parameter must be sufficiently large to ensure equivalence, allowing the input embeddings to dominate the prior, but not so large that it significantly prolongs the regularisation. Previously, it was set as a ratio of the empirical standard deviation from training data (Fehr & Henderson, 2024). In this work, we determine it manually by selecting the smallest  $\tau_\alpha$  that maintains equivalence of the pretrained model.

$$\boldsymbol{\alpha} = \alpha(\mathbf{x}) = \exp(\mathbf{x}^2\mathbf{w}_1^\alpha + \mathbf{x}\mathbf{w}_2^\alpha + b^\alpha); \quad \mathbf{w}_1^\alpha = \frac{1}{2\sqrt{d/h}} \odot \mathbf{1}; \quad \mathbf{w}_2^\alpha = \mathbf{0}; \quad b^\alpha = \tau_\alpha \quad (19)$$

**Fine-tuning hyperparameters.** Following the approach of Henderson & Fehr (2023), we set the number of samples per component to  $\kappa^\Delta = 1$ . However, the authors define a conditional prior, which when training models from scratch helped to control the sparsity. In this work, we do not incorporate this conditional prior. As shown in Appendix A.3, the Kullback-Leibler divergence is decomposed into two parts, with the Gaussian and Dirichlet components weighted by the hyperparameters  $\lambda_G$  and  $\lambda_D$ , respectively. We explore different of these hyperparameters during fine-tuning for each experiment.

### C.1 SPEECH OUT-OF-DISTRIBUTION EVALUATION

#### C.1.1 SPEECH QUALITY ASSESSMENT

**Fine-tuning details.** For the speech quality regression task on NISQA<sup>2</sup> and Tencent<sup>3</sup>, we used the mean-squared-error (MSE) loss with the pretrained Wav2vec2-base<sup>4</sup> model (Baeovski et al., 2020), a 12 Transformer encoder. For the regression head we use two linear layers, including non-linearity and dropout, followed by mean pooling. The size for the latent embedding vectors and model projections are 768 with 12 attention heads, which leads to models of approximately 95 million parameters. During fine-tuning we use: the Adam optimiser (Kingma & Ba, 2014), a constant learning rate of  $1e^{-5}$ , batch size of 16, trained for 5 epochs. For regularised models we include layer drop of 0.1 time-frequency masking on the output of the feature encoder with probability 0.05. The pretrained convolutional feature encoder is not fine-tuned. This is a standard for the model architecture (Baeovski et al., 2020).

<sup>2</sup>Dataset: <https://github.com/gabrielmittag/NISQA/wiki/NISQA-Corpus>

<sup>3</sup>Dataset: <https://github.com/ConferencingSpeech/ConferencingSpeech2022>

<sup>4</sup>Model: <https://huggingface.co/facebook/wav2vec2-base-960h>

**NVIB details.** For our Transformer encoder, we include NVIB in layers 0–10. The NVIB projections were initialised using  $\tau_\sigma = 0.1$  and  $\tau_\alpha = 10$ . During fine-tuning we included the learnable prior  $\mu^p$ . Hyperparameters for influencing the amount of regularisation from the Gaussian component  $\lambda_G$  and Dirichlet component  $\lambda_D$  were tied and selected over a log-scaled grid search  $[1e^{-1}, 1e^{-2}, 1e^{-3}, 1e^{-4}]$  on validation data. The best performing model used NVIB regularisation parameters of  $\lambda_G = \lambda_D = 1e^{-2}$  and evaluated over 5 seeds.

### C.1.2 SPEECH LANGUAGE IDENTIFICATION

**Fine-tuning details.** For the language identification classification task, we used cross-entropy loss with the pretrained Wav2vec2-large<sup>5</sup> model (Baevski et al., 2020), a 24 layer Transformer encoder. We use a single linear layer and mean pooling classification head. The size for the latent embedding vectors and model projections are 1024 with 16 attention heads, which leads to models of approximately 317 million parameters. During fine-tuning we use: the AdamW (Loshchilov & Hutter, 2019) optimiser, a learning rate of  $3e^{-5}$ , scheduler with a linear warm-up and decay, batch size of 4, trained for 10 epochs with mixed precision 16bit and gradient norm clipping of 1. This experiment also includes weight decay for all models of 0.05. For regularised models, we include layer drop of 0.1 and time-frequency masking on the output of the feature encoder with probability 0.05. The pretrained convolutional feature encoder is not fine-tuned. This is a standard for the model architecture (Baevski et al., 2020).

**NVIB details.** For our Transformer encoder, we include NVIB in layers 0–16. The NVIB projections were initialised using  $\tau_\sigma = 0$  and  $\tau_\alpha = 10$ . During fine-tuning we included the learnable prior  $\mu^p$ . Hyperparameters for influencing the amount of regularisation from the Gaussian component  $\lambda_G$  and Dirichlet component  $\lambda_D$  were tied and selected over a log-scaled grid search  $[1e^{-1}, 1e^{-2}, 1e^{-3}, 1e^{-4}, 1e^{-5}, 1e^{-6}, 1e^{-7}, 1e^{-8}]$  on validation data. The best performing model used NVIB regularisation parameters of  $\lambda_G = \lambda_D = 1e^{-7}$  and evaluated over 5 seeds.

**Dataset Details.** We fine-tune our models on the CommonLanguage<sup>6</sup> (Ravanelli et al., 2021) dataset which consists of 22K training audios from 45 languages. For VoxPopuli<sup>7</sup> we selected 11 languages (Czech, Dutch, English, Estonian, French, German, Italian, Polish, Romanian, Slovenian, and Spanish) while 27 for FLEURS<sup>8</sup> (Arabic, Catalan, Czech, Dutch, English, Estonian, French, Georgian, German, Greek, Indonesian, Italian, Japanese, Kyrgyz, Latvian, Maltese, Persian, Polish, Portuguese, Romanian, Russian, Slovenian, Spanish, Swedish, Tamil, Turkish, and Welsh).

### C.2 TEXT OUT-OF-DISTRIBUTION CLASSIFICATION

**Fine-tuning details.** For the CivilComments<sup>9</sup> classification task, we used cross-entropy loss with the pretrained TinyBERT<sup>10</sup> model (Turc et al., 2019), a two-layer Transformer encoder. The size for the latent embedding vectors and model projections are 128 with 2 attention heads, which leads to models of approximately 4.5 million trainable parameters. The standard BERT base-uncased tokenizer is used for tokenisation with a vocabulary of approximately 30K. During fine-tuning we use: the AdamW optimiser (Loshchilov & Hutter, 2019), a constant learning rate of  $5e^{-5}$ , batch size of 1024, trained for 50 epochs with mixed precision 16bit and gradient norm clipping of 0.1.

**NVIB details.** For our two-layer Transformer encoder, we include NVIB in both layers. The NVIB projections were initialised using  $\tau_\sigma = 0.1$  and  $\tau_\alpha = 1$ . During fine-tuning a linear KL annealing warmup was used including the learnable prior  $\mu^p$ . Hyperparameters for influencing the amount of regularisation from the Gaussian component  $\lambda_G$  and Dirichlet component  $\lambda_D$  were tied and selected over a log-scaled grid search  $[1e^0, 1e^{-1}, 1e^{-2}, 1e^{-3}, 1e^{-4}, 1e^{-5}]$  on validation data. The best performing model used NVIB regularisation parameters of  $\lambda_G = \lambda_D = 1e^{-1}$  and evaluated over 5 seeds.

**Induced sparsity.** After fine-tuning, NVIB regularisation naturally induces sparsity in the key-space, whereas dropout promotes a more uniform distribution across keys, as shown in the right-most frames of the attention maps (Figures 4 and 5). This effect arises from the NVIB regularisation decreasing the weight of embeddings in proportion to the prior component embedding during the attention calculation. To remove a key, we first calculate the average attention weights across the query dimension and then mask out the embeddings with the lowest magnitudes, thereby inducing key sparsity. Since NVIB naturally creates key sparsity, when

<sup>5</sup>Model: <https://huggingface.co/facebook/wav2vec2-large-960h>

<sup>6</sup>Dataset: [https://huggingface.co/datasets/speechbrain/common\\_language](https://huggingface.co/datasets/speechbrain/common_language)

<sup>7</sup>Dataset: <https://huggingface.co/datasets/facebook/voxpathuli>

<sup>8</sup>Dataset: <https://huggingface.co/datasets/google/fleurs>

<sup>9</sup>Dataset: <https://github.com/p-lambda/wilds>

<sup>10</sup>Model: [https://huggingface.co/google/bert\\_uncased\\_L-2\\_H-128\\_A-2](https://huggingface.co/google/bert_uncased_L-2_H-128_A-2)

these keys are dropped, we notice an improvement or sustained task performance (Figure 2). As the proportion of keys being masked increases (right to left in Figures 4 and 5), we notice minor changes in the attention patterns for the NVIB model and clear alignment with the tokens that are toxic content.

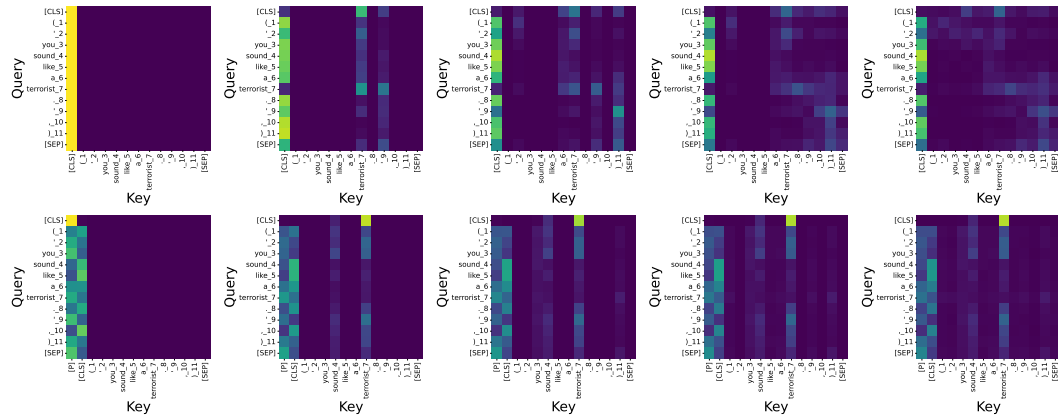


Figure 4: Attention plot for the best models on CivilComments. The plots show a single head of the last layer. Left-Right: Proportion of keys retained [0.1, 0.25, 0.5, 0.75, 1.0]. Top: with Dropout. Bottom: with NVIB. Sentence: (‘you sound like a terrorist.’). NVIB highlights ‘sound’ and ‘terrorist’.

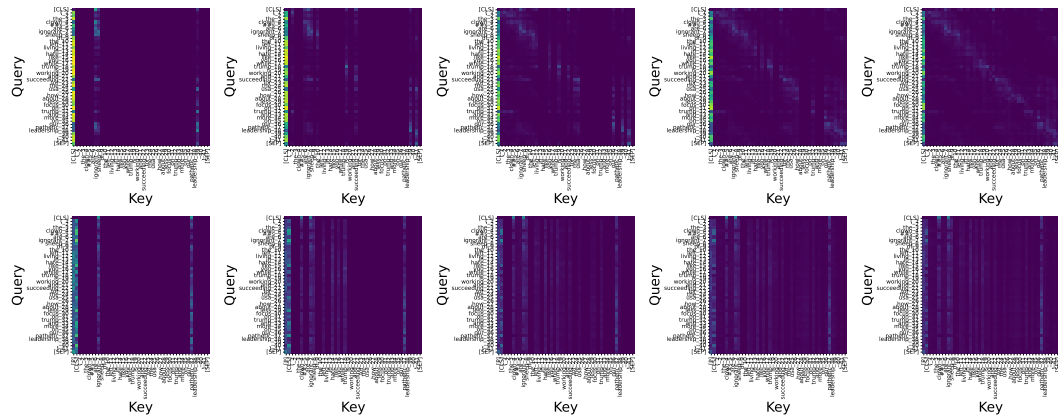


Figure 5: Attention plot for the best models on CivilComments. The plots show a single head of the last layer. Left-Right: Proportion of keys retained [0.1, 0.25, 0.5, 0.75, 1.0]. Top: with dropout. Bottom: with NVIB. Sentence: (‘the clowns are ignorant sheep of the left living in hate like you while trump is working and succeeding for the usa. how about less focus on trump and more on our pathetic leadership’). NVIB highlights ‘ignorant’ and ‘pathetic’.

### C.3 GRAPH LINK PREDICTION

We build upon the BERT for Link Prediction (BLP) model (Daza et al., 2020), which encodes the textual description of  $h$  and  $t$  using BERT. The resulting representations are pooled into the final layer’s [CLS] dense representation, denoted as  $h_{[\text{CLS}]}$  and  $t_{[\text{CLS}]}$ , respectively. The relation  $r$  is selected from a learnable lookup table  $R$ . The model is trained using a contrastive approach, where a positive triple  $(h, r, t)$  is compared to corrupted triples  $(h', r, t')$  using a distance-based loss function, such as TransE (Bordes et al., 2013):

$$f_{\text{TransE}}(h, r, t) = ||h + r - t||$$

We adopt the experimental setting of BLP without further hyperparameter tuning.

**Fine-tuning details.** For the graph-link classification task on FB15k-237<sup>11</sup>, we used distance-based loss function model, with the pretrained TinyBERT<sup>12</sup> model (Turc et al., 2019), a two-layer Transformer encoder. The size for the latent embedding vectors and model projections are 128 with 2 attention heads, which leads to models of approximately 4.5 million trainable parameters. The standard BERT base-uncased tokenizer is used for tokenisation with a vocabulary of approximately 30K. During fine-tuning we use: the RAdam optimiser (Liu et al., 2020), a cosine learning rate scheduler with value  $8e^{-5}$ , batch size of 256, trained for 40 epochs with mixed precision 16bit and gradient norm clipping of 1.

**NVIB details.** For our two-layer Transformer encoder, we include NVIB in both layers. The NVIB projections were initialised using  $\tau_\sigma = 0.1$  and  $\tau_\alpha = 1$ . During fine-tuning the learnable prior  $\mu^p$  was used. Hyperparameters for influencing the amount of regularisation from the Gaussian component  $\lambda_G$  and Dirichlet component  $\lambda_D$  were tied and selected over a log-scaled grid search  $[1e^{-1}, 1e^{-2}, 1e^{-3}, 1e^{-4}]$  on validation data. The best performing model used NVIB regularisation parameters of  $\lambda_G = \lambda_D = 1e^{-3}$ .

### C.4 IMAGE FEW SHOT CLASSIFICATION

These experiments build from the following work and repository Hu et al. (2022)<sup>13</sup>.

#### C.4.1 FEW-SHOT IN-DISTRIBUTION

**Fine-tuning details.** For the CIFAR-FS<sup>14</sup> few-shot classification task, we used cross-entropy loss with the pretrained DeiT-Small<sup>15</sup> model (Touvron et al., 2021a), a 12 Transformer encoder. The size for the latent embedding vectors and model projections are 384 with 6 attention heads, which leads to models of approximately 22 million trainable parameters. We use the prototypical network (ProtoNet) (Snell et al., 2017) classifier, which creates class centroids dynamically for each episode and then performs nearest centroid classification (Hu et al., 2022). During fine-tuning we use: the AdamW optimiser (Loshchilov & Hutter, 2019), a linear warmup with cosine decay learning rate scheduler  $1e^{-4}$ , batch size of 1, trained for 50 epochs with mixed 16bit precision. Episodes during meta-training 2000 and episodes during testing is 2000. Weight decay is kept constant 0.05 for all experiments.

**NVIB details.** For our Transformer encoder, we include NVIB in layers 0 – 5. The NVIB projections were initialised using  $\tau_\sigma = 0$  and  $\tau_\alpha = 0$ . In this experiment we initialised the prior  $\mu^p = \mathbf{0}$  and did not allow it to be learnable. Hyperparameters for influencing the amount of regularisation from the Gaussian component  $\lambda_G$  and Dirichlet component  $\lambda_D$  were tied and selected over a log-scaled grid search  $[1e^0, 1e^{-1}, 1e^{-2}, 1e^{-3}, 1e^{-4}, 1e^{-5}]$  on validation data. The best performing model used NVIB regularisation parameters of  $\lambda_G = \lambda_D = 1e^{-2}$ .

#### C.4.2 FEW-SHOT OUT-OF-DISTRIBUTION

This experiment uses the Meta-Dataset (Triantafillou et al., 2019), which samples 5 – 50 ways, with a maximum support size of 500 and a maximum query size of 10. For datasets except ImageNet and Omniglot, we use uniform sampling, while for ImageNet and Omniglot, sampling is performed according to the hierarchy of classes. This follows the methodology described in Hu et al. (2022).

<sup>11</sup>Dataset: <https://github.com/dfdazac/blp?tab=readme-ov-file>

<sup>12</sup>Model: [https://huggingface.co/google/bert\\_uncased\\_L-2\\_H-128\\_A-2](https://huggingface.co/google/bert_uncased_L-2_H-128_A-2)

<sup>13</sup>Original codebase: [https://github.com/hushell/pmf\\_cvpr22](https://github.com/hushell/pmf_cvpr22)

<sup>14</sup>Dataset dropbox link: <https://www.dropbox.com/scl/fi/91dgxywb8e948rvmmq1d8/cifar-fs-splits.zip?rlkey=h69z5fxhelrdonjm9s37q6laf&e=1&dl=0>

<sup>15</sup>Model: <https://huggingface.co/facebook/deit-small-patch16-224>

**Fine-tuning details.** For the Meta-Dataset<sup>16</sup> few-shot classification task, we used cross-entropy loss with the pretrained DeiT-Small<sup>17</sup> model (Touvron et al., 2021a), a 12 Transformer encoder. The size for the latent embedding vectors and model projections are 384 with 6 attention heads, which leads to models of approximately 22 million trainable parameters. We use the prototypical network (ProtoNet) (Snell et al., 2017) classifier, which creates class centroids dynamically for each episode and then performs nearest centroid classification (Hu et al., 2022). During fine-tuning we use: the AdamW optimiser (Loshchilov & Hutter, 2019), a linear warmup with cosine decay learning rate scheduler  $1e^{-4}$ , batch size of 1, trained for 50 epochs with mixed 16bit. Episodes during meta-training 2000 and episodes during testing is 600. Weight decay is kept constant 0.05 for all experiments.

**NVIB details.** For our Transformer encoder, we include NVIB in layers 0 – 5. The NVIB projections were initialised using  $\tau_\sigma = 0$  and  $\tau_\alpha = -3$ . In this experiment we initialised the prior  $\mu^p = \mathbf{0}$  and did not allow it to be learnable. Hyperparameters for influencing the amount of regularisation from the Gaussian component  $\lambda_G$  and Dirichlet component  $\lambda_D$  were tied and selected over a log-scaled grid search  $[1e^0, 1e^{-1}, 1e^{-2}, 1e^{-3}, 1e^{-4}, 1e^{-5}]$  on validation data. The best performing model used NVIB regularisation parameters of  $\lambda_G = \lambda_D = 1e^{-3}$ .

---

<sup>16</sup>Dataset: <https://github.com/google-research/meta-dataset>

<sup>17</sup>Model: <https://huggingface.co/facebook/deit-small-patch16-224>