

EFFICIENT UNCERTAINTY ESTIMATION WITH GAUSSIAN PROCESS FOR RELIABLE DIALOG RESPONSE RETRIEVAL

Tong Ye^{1,2}, Zhitao Li¹, Jianzong Wang^{1,*}, Ning Cheng¹, Jing Xiao¹

¹Ping An Technology (Shenzhen) Co., Ltd.

²University of Science and Technology of China

ABSTRACT

Deep neural networks have achieved remarkable performance in retrieval-based dialogue systems, but they are shown to be ill calibrated. Though basic calibration methods like Monte Carlo Dropout and Ensemble can calibrate well, these methods are time-consuming in the training or inference stages. To tackle these challenges, we propose an efficient uncertainty calibration framework GPF-BERT for BERT-based conversational search, which employs a Gaussian Process layer and the focal loss on top of the BERT architecture to achieve a high-quality neural ranker. Extensive experiments are conducted to verify the effectiveness of our method. In comparison with basic calibration methods, GPF-BERT achieves the lowest empirical calibration error (ECE) in three in-domain datasets and the distributional shift tasks, while yielding the highest $R_{10}@1$ and MAP performance on most cases. In terms of time consumption, our GPF-BERT has an $8\times$ speedup.

Index Terms— Uncertainty, Calibration, Gaussian Process, Dialog Response Retrieval

1. INTRODUCTION

Dialog response retrieval models based on deep neural networks (DNNs) primarily focus on modeling the relevance between context and responses and have achieved impressive performance [1, 2, 3]. However, these models always suffered from over or under confidence due to the poor calibration of DNNs [4]. As a result, it is difficult to determine whether the predictions are reliable. This attribute is essential for distribution shifts tasks and safety-critical areas since erroneous predictions can result in far more significant consequences than not making any prediction at all [5]. Therefore, an ideal dialog model should exhibit confidence in its predictions while also recognizing situations where its predictions may be incorrect and uncertain.

Uncertainty modeling has been touched in previous work on dialog systems. Monte Carlo (MC) Dropout [2, 3, 6] and Ensemble [2, 7] have emerged as two of the most prominent uncertainty estimation methods for deep retrieval networks. While Ensemble trains independently multiple models

using stochastic gradient descent, MC Dropout trains a single stochastic network by dropping different subsets of weights simultaneously in train and test time [7]. Unfortunately, MC Dropout necessitates carrying out several forward passes and Ensemble becomes computationally expensive. This poses a significant challenge, particularly in light of the widespread adoption of large transformer architectures like BERT [7]. Therefore, it is urgent to explore an efficient method to quantify the uncertainty in deep neural retrieval models.

Gaussian Process (GP) [8] is a well-established framework for evaluating uncertainty. As an input moves farther away from the training data, the level of uncertainty in GP predictions tends to increase [9]. However, GP is challenging to scale to large datasets and improve the performance while DNNs are computationally scalable enough to handle them [9]. SNGP [10] combines the strengths of GP and DNNs, utilizing spectral normalization [11] to the weights in each residual layer, which can efficiently handle the large scale inputs and makes robust uncertainty-aware predictions. Unfortunately, the application of SNGP has not been explored in the retrieval-based dialog system.

So motivated, we attempt to investigate a simple and efficient approach for the well-calibrated dialog response retrieval models based on Gaussian Process. Specifically, we add a neural GP layer to a deterministic BERT-like backbone to improve the ability of uncertainty estimation and train the model with focal loss [12] to achieve better calibration. Different from MC Dropout and Ensemble, our method only needs to be performed by passing through a single forward so that GPF-BERT achieves almost $8\times$ speedup in terms of inference time. To summarize, the main contributions of this work are as follows:

- To our best knowledge, we first estimate uncertainty in dialog tasks with SNGP. Furthermore, we propose an efficient framework GPF-BERT to estimate uncertainty combining the focal loss and SNGP.
- We conduct extensive experiments to compare the performance of various calibration methods. Our method yields the lowest ECE in three in-domain datasets and the distributional shift task while keeping performance.

*Corresponding author: Jianzong Wang, jzwang@188.com.

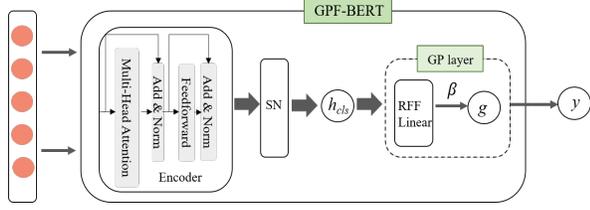


Fig. 1. An illustration of GPF-BERT prediction models for dialog response retrieval.

2. METHODOLOGY

In this section, we present details of GPF-BERT, utilizing a neural Gaussian Process layer to model uncertainty as shown in Fig 1.

2.1. Model Architecture

For a dialog dataset, we denote a training set as a triples $\{(U_i, r_i, y_i)\}_{i=1}^N$, where $U_i = \{u_{i1}, u_{i2}, \dots, u_{it}\}$ is a dialog context consisting of t utterances and y_i is the response relevance label $y_i \in \{0, 1\}$, with a response candidate r_i . The whole dialog context with the candidate is fed into BERT-like encoders, $x_i = \{[CLS], u_{i1}, \dots, u_{it}, [SEP], r_i\}$, where the special token $[CLS]$ denotes the sequence beginning and $[SEP]$ separates the response from the contexts. For each training example, the representation vector $h(x_i)$ of $[CLS]$ is the feature of this dialog.

For a latent representation $h_i = h(x_i)$, the GP output layer $g_i = g(h_i)$ follows a multivariate normal distribution a priori: $g \sim GP(0, K)$, where K is a $N \times N$ kernel matrix. For more details of GPs, refer to [8].

First, spectral normalization (SN) [11] involves decomposing the parameters W of each neural network layer using SVD and subsequently constraining the maximum singular value to 1. Briefly, the SN estimates the spectral norm $\|W_l\|_2$ using the power iteration method during each training step and subsequently normalizes the weights according to the estimated norm as follows:

$$w_l = \begin{cases} c * W_l / \|W_l\|_2 & \text{if } c < \|W_l\|_2 \\ W_l & \text{otherwise} \end{cases} \quad (1)$$

To approximate the kernel matrix with a low-rank approximation, GPF-BERT uses a technique known as random Fourier features (RFF) [13] as $K = \Phi\Phi^T$ ($\Phi \in R^{N \times L}$), where L denotes the dimensionality of the latent space:

$$g \sim GP(0, \Phi\Phi^T), \Phi_i = \sqrt{2/L} * \cos(-Wh_i + b) \quad (2)$$

Φ_i contains a fixed weight matrix W with entries sampled from $N(0, 1)$, and a fixed bias vector b with entries sampled from $U(0, 2\pi)$.

Finally, the GP layer as a neural network layer [14] with learnable output weights $\beta \sim N(0, I)$ according to the RFF

approximation to the GP prior.

$$g(h_i) = \sqrt{2/L} * \cos(-Wh_i + b)^T \beta \quad (3)$$

Note that the use of random Fourier features (RFF) to linearly transform the GP helps to overcome the curse of dimensionality and yield good performance with finite data. However, this technique may also result in feature collapse, which can compromise model robustness, particularly when dealing with distributional shifts. Therefore, we deploy SN on the weights, which enforces the feature extractor to be bi-Lipschitz to mitigate feature collapse [15].

Due to the lack of conjugacy between the classified likelihood function and a Gaussian prior, we resort to using the Laplace approximation to quantify uncertainty [10]. The maximum a posteriori (MAP) solution is denoted by $\hat{\beta}$ and the Laplace posterior for GP under the RFF approximation can be expressed as:

$$p(\beta|D) \sim GP(\hat{\beta}, \hat{\Sigma}_k), \hat{\Sigma}_k^{-1} = \sum_{i=1}^N \hat{p}_i(1 - \hat{p}_i)\Phi_i\Phi_i^T + I \quad (4)$$

During minibatch training, $\hat{\beta}$ is updated via regular SGD with respect to the loss function and $\hat{\Sigma}_t^{-1}$ is updated cheaply using

$$\hat{\Sigma}_t^{-1} = \alpha * \hat{\Sigma}_{t-1}^{-1} + (1 - \alpha) * \sum_{i=1}^M \hat{p}_i(1 - \hat{p}_i)\Phi_i\Phi_i^T \quad (5)$$

where t indexes update steps, M is the mini-batch size, \hat{p}_i is the softmax probability and α is a small scaling coefficient.

For a given feature vector x^* of a query-response pair, GPF-BERT computes the posterior mean $\hat{m}(x^*) = \Phi^T \hat{\beta}$ and the variance $\hat{K}(x^*) = \Phi^T \hat{\Sigma} \Phi$. Finally, the predictive distribution is written as $p = \exp(m) / \sum_i \exp(m_i)$ where $m \sim N(\hat{m}, \hat{K})$ and we calculate its posterior mean using mean-field approximation [16] for lower computation.

2.2. Loss Function

The conventional cross-entropy loss assigns the same weight to individual samples on a mini-batch. However, there are several low-confidence samples hard to classify. Although the high-confidence samples have a small loss, their cumulative loss value is still greater than the low-confidence samples due to the large number, dominating gradient and producing bad performance. That is just an important reason for poor calibration of cross-entropy [17].

In this work, our GPF-BERT utilizes focal loss [12], which focuses more on the uncertain samples by reducing the weight of high-confidence samples. Namely, by reducing the weight of the easy samples, the model focuses more on the hard samples when training.

$$\begin{aligned} L_{focal} &= -(1 - \hat{p})^\gamma \log(\hat{p}) \geq -(1 - \gamma\hat{p}) \log \hat{p} \\ &= L_{ce} - \gamma H[\hat{p}] \geq KL(q||\hat{p}) - \gamma H[\hat{p}] \end{aligned} \quad (6)$$

When training, focal loss ensures minimization of the KL divergence whilst simultaneously increasing the entropy $H[\hat{p}]$

[17]. The high entropy can help prevent the model from becoming overconfident and thereby improve calibration. We use L_{ce} and L_{focal} to respectively denote the cross-entropy loss and the focal loss with hyperparameter $\gamma \geq 1$.

3. EXPERIMENTS

3.1. Experiments Setup

Datasets: We utilize three large-scale conversational response ranking datasets in our experiments. MS Dialog [18] contains 246,000 context-response pairs culled from over the Microsoft Answer community. MANTIS [19] contains 1.3 million context-response pairs including 14 different domain. Ubuntu Dialogue Corpus v1.0 (UDC) [20] is consisted of almost 1 million context-response pairs.

Metrics: We use $R_{10}@1$ and MAP to measure the retrieval performance and ECE [21] for calibration. We divide the interval $[0, 1]$ into $M = 10$ equispaced bins. The ECE calculate a weighted average of the absolute difference between the accuracy A_i and confidence B_i of each bin: $ECE = \sum_{i=1}^M \frac{|B_i|}{N} |mean(A_i) - mean(B_i)|$.

Baselines: We compare BERT [1] with other common calibrated methods. MC Dropout [2] approximates Bayesian inference using dropout during training and testing, and produces a predictive distribution by performing multiple forward passes. In this paper, we use 10 forward passes based on a BERT model. Ensemble [2] independently trained several models and integrate their predictions a model integrating predictions of BERT and MC Dropout models. SNGP [10] adds a weighting normalization step during training and replaces the dense output layer with a GP layer.

Implementation Details: We use 12-layered *BERT* as the backbone, each encoder having 12 attention heads and a hidden dimension of 768. All the methods are obtained a latent representation by extracting the $[CLS]$ feature. The dropout probability was set to 0.1 and the learning rate was set to $5e-6$ using the Adam optimizer. Following recent research [2] that employed finetuned BERT for dialog response ranking, we randomly select nine responses from the list of all responses as the negative samples when training. The hyperparameters utilized for the deterministic variant are also employed for each model architecture. In our experiment, we first train our model for 1 epochs with a batch size of 16 on a cluster of 1 Tesla V100 with 16G memory. Simultaneously, all the experiment results are the average over 5 runs along with the standard error.

3.2. Results

In-domain. Table 1 reports the results of GPF-BERT and all baselines on the MS dialog, MANTIS and UDC datasets. From this table, GPF-BERT generally outperforms other

Table 1. Calibration (ECE) and effectiveness ($R_{10}@1$, MAP). ” \uparrow ” represents higher is better and ” \downarrow ” means lower is better.

		$R_{10}@1\uparrow$	MAP \uparrow	ECE \downarrow
MSDialog	BERT	0.682\pm0.006	0.800 \pm 0.003	0.125 \pm 0.020
	MC Dropout	0.673 \pm 0.005	0.796 \pm 0.003	0.110 \pm 0.020
	Ensemble	0.680 \pm 0.004	0.800\pm0.003	0.115 \pm 0.019
	SNGP	0.659 \pm 0.013	0.783 \pm 0.008	0.110 \pm 0.006
	GPF-BERT	0.681 \pm 0.006	0.799 \pm 0.003	0.025\pm0.010
MANTIS	BERT	0.590 \pm 0.012	0.713 \pm 0.010	0.169 \pm 0.028
	MC Dropout	0.591 \pm 0.011	0.713 \pm 0.009	0.152 \pm 0.02
	Ensemble	0.592 \pm 0.011	0.713 \pm 0.010	0.157 \pm 0.026
	SNGP	0.597 \pm 0.021	0.719 \pm 0.013	0.147 \pm 0.010
	GPF-BERT	0.614\pm0.017	0.729\pm0.012	0.025\pm0.009
UDC	BERT	0.810 \pm 0.003	0.880 \pm 0.002	0.037 \pm 0.001
	MC Dropout	0.809 \pm 0.003	0.878 \pm 0.002	0.033 \pm 0.001
	Ensemble	0.810 \pm 0.003	0.879 \pm 0.002	0.034 \pm 0.001
	SNGP	0.806 \pm 0.001	0.877 \pm 0.001	0.033 \pm 0.001
	GPF-BERT	0.818\pm0.001	0.885\pm0.001	0.016\pm0.002

single-model approaches in ECE across various datasets with a reduction of less than 1% in $R_{10}@1$ and MAP.

Specifically, BERT, which is a vanilla model without any calibration, usually performs well but is not well calibrated. The calibration of MC Dropout and Ensemble exceeds BERT, which verifies that Bayesian models exhibit greater expressiveness in their ability to convey confidence, but unfortunately still obtain poor calibration. Compared to the prior methods, the ECE of GPF-BERT is almost lowest, which is reduced by almost 10%, 14% and 2% respectively in three datasets, while the $R_{10}@1$ and MAP are better or less than 1% decrease. Namely, GPF-BERT includes uncertainty information while keeping $R_{10}@1$ and MAP performance in in-domain datasets.

Distributional Shift. In addition to using the uncertainty estimation for in-domain datasets, we also train the model using the training set from one dataset, i.e. train set, and evaluate it on a different dataset’s test set, which is also known as domain generalization or distributional shift tasks. We record all the retrieval performance and calibration in Table 2. As shown, we observe that SNGP was reduced by up to 4% when compared to the BERT, MC Dropout and Ensemble. Moreover, the GPF-BERT achieves a substantial decrease up to 10% to the upper calibration bound under the SNGP framework, even though the $R_{10}@1$ and MAP of GPF-BERT is lower in distribution shift tasks. This confirms that GP-based retrieval models will have the more robust expressiveness to convey confidence in distributional shift tasks. According to the results of SNGP and GPF-BERT, we find that the focal loss plays an important role in calibration.

Efficiency. One of the most critical challenges to overcome is the computational cost when employing Bayesian to capture uncertainty. We analyze the efficiency of GPF-BERT in terms of parameter number and inference time in the MS Dialog dataset in Table 3. Compared to MC Dropout and Ensem-

Table 2. Calibration (ECE) and effectiveness ($R_{10}@1$, MAP) for distributional shift tasks. ” \uparrow ” represents higher is better and ” \downarrow ” means lower is better. All the models are trained in one dataset and test in the other dataset.

Train	Test	Metric	BERT	MC Dropout	Ensemble	SNGP	GPF-BERT
MS Dialog	MANtIS	$R_{10}@1\uparrow$	0.378 \pm 0.024	0.357 \pm 0.017	0.369 \pm 0.020 \pm	0.385\pm0.040	0.381 \pm 0.033
		MAP \uparrow	0.538 \pm 0.018	0.524 \pm 0.012	0.533 \pm 0.015	0.543\pm0.030	0.540 \pm 0.023
		ECE \downarrow	0.343 \pm 0.035	0.328 \pm 0.045	0.331 \pm 0.037	0.307 \pm 0.011	0.206\pm0.054
	UDC	$R_{10}@1\uparrow$	0.609 \pm 0.009	0.600 \pm 0.007	0.606 \pm 0.008	0.602 \pm 0.007	0.612\pm0.014
		MAP \uparrow	0.736 \pm 0.005	0.730 \pm 0.004	0.734 \pm 0.004	0.731 \pm 0.004	0.737\pm0.009
		ECE \downarrow	0.109 \pm 0.011	0.092 \pm 0.011	0.097 \pm 0.009	0.085 \pm 0.003	0.026\pm0.016
MANtIS	MS Dialog	$R_{10}@1\uparrow$	0.430 \pm 0.069	0.418 \pm 0.060	0.427 \pm 0.065	0.363 \pm 0.108	0.452\pm0.037
		MAP \uparrow	0.598 \pm 0.058	0.591 \pm 0.050	0.596 \pm 0.054	0.548 \pm 0.089	0.619\pm0.026
		ECE \downarrow	0.514 \pm 0.037	0.497 \pm 0.040	0.503 \pm 0.037	0.485 \pm 0.017	0.364\pm0.019
	UDC	$R_{10}@1\uparrow$	0.662 \pm 0.008	0.660 \pm 0.010	0.662\pm0.009	0.656 \pm 0.013	0.659 \pm 0.015
		MAP \uparrow	0.769 \pm 0.008	0.768 \pm 0.006	0.769\pm0.006	0.766 \pm 0.009	0.768 \pm 0.009
		ECE \downarrow	0.071 \pm 0.006	0.061 \pm 0.006	0.064 \pm 0.006	0.059 \pm 0.008	0.025\pm0.011
UDC	MS Dialog	$R_{10}@1\uparrow$	0.324 \pm 0.031	0.290 \pm 0.065	0.312 \pm 0.072	0.406 \pm 0.092	0.431\pm0.074
		MAP \uparrow	0.513 \pm 0.072	0.484 \pm 0.064	0.503 \pm 0.069	0.581 \pm 0.072	0.603\pm0.059
		ECE \downarrow	0.616 \pm 0.031	0.622 \pm 0.026	0.610 \pm 0.030	0.607 \pm 0.034	0.528\pm0.042
	MANtIS	$R_{10}@1\uparrow$	0.250 \pm 0.045	0.220 \pm 0.038	0.240 \pm 0.044	0.270 \pm 0.027	0.283\pm0.035
		MAP \uparrow	0.418 \pm 0.037	0.389 \pm 0.032	0.408 \pm 0.036	0.437 \pm 0.021	0.447\pm0.029
		ECE \downarrow	0.510 \pm 0.041	0.537 \pm 0.037	0.508 \pm 0.040	0.494 \pm 0.032	0.400\pm0.045

Table 3. Parameters and inference time.

models	BERT	MC Dropout	Ensemble	GPF-BERT
Parameters(M)	413.26	413.26	826.52	453.55
Time(min)	6.80(1 \times)	65.00(9.56 \times)	71.80(10.56 \times)	7.87(1.16 \times)

ble, GPF-BERT has a significant decrease (at least 8 times) in inference time. While not completely free, GPF-BERT only adds negligible computational cost, but it greatly improves the calibration, which facilitates adaptation to other models.

We believe that GP maintains a distribution over functions rather than model parameters, which enables GPF-BERT to improve uncertainty calibration for dialog response retrieval models. In addition, according to a recent theorem [22] that capturing uncertainty information and correcting overconfidence can be achieved by making only the last layer of a model in binary classification, we can assume that adding a GP layer is Bayesian enough so that GPF-BERT can achieve better calibration.

3.3. Ablation Study

To understand the impact of focal loss on calibration improvement, we conducted a straightforward ablation study on the in-domain dataset MS Dialog and its two distribution shift tasks as shown in Fig 2. Obviously, the $R_{10}@1$ and MAP of GPF-BERT are lower than Focal but the calibration is much better. Focal loss is a commonly used regularization method and we found that it still facilitates calibration when applied to SNGP. On the other hand, the difference in architecture between Focal and GPF-BERT demonstrates that the GP-based model architecture is a better framework for calibration than conventional models. That is to say, the model architecture may be key to improving calibration.

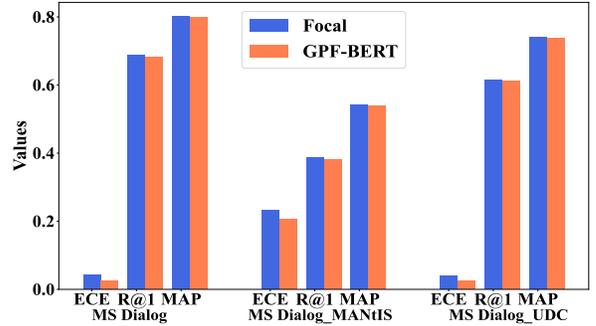


Fig. 2. Ablation study.

4. CONCLUSION

In this paper, we present an efficient uncertainty estimation architecture GPF-BERT for reliable dialog response retrieval tasks. GPF-BERT only adds a neural GP layer to a deterministic DNN to improve the ability of uncertainty estimation and trains the model with focal loss to achieve better calibration while maintaining the flexibility of deep neural networks. We conducted extensive experiments to verify the effectiveness including parameters and inference time. Furthermore, we explored the relative contributions of focal loss to the effectiveness improvement in the ablation study.

5. ACKNOWLEDGEMENT

This paper is supported by the Key Research and Development Program of Guangdong Province under grant No. 2021B0101400003. Jianzong Wang from Ping An Technology (Shenzhen) Co., Ltd (jzwang@188.com) is the corresponding author.

6. REFERENCES

- [1] 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*, 2019, pp. 4171–4186.
- [2] Gustavo Penha and Claudia Hauff, “On the calibration and uncertainty of neural learning to rank models for conversational search,” in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, 2021, pp. 160–170.
- [3] Daniel Cohen, Bhaskar Mitra, Oleg Lesota, Navid Reksabsaz, and Carsten Eickhoff, “Not all relevance scores are equal: Efficient uncertainty and calibration modeling for deep retrieval models,” in *Proceedings of the 44th International ACM SIGIR Conference*, 2021, pp. 654–664.
- [4] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger, “On calibration of modern neural networks,” in *ICML*, 2017, pp. 1321–1330.
- [5] Amita Kamath, Robin Jia, and Percy Liang, “Selective question answering under domain shift,” in *ACL, July 5-10, 2020*, pp. 5684–5696.
- [6] Yarin Gal and Zoubin Ghahramani, “Dropout as a bayesian approximation: Representing model uncertainty in deep learning,” in *ICML*, 2016, pp. 1050–1059.
- [7] Nikita Durasov, Timur Bagautdinov, Pierre Baque, and Pascal Fua, “Masksembles for uncertainty estimation,” in *CVPR*, 2021, pp. 13539–13548.
- [8] Christopher K Williams and Carl Edward Rasmussen, *Gaussian processes for machine learning*, vol. 2, MIT press Cambridge, MA, 2006.
- [9] Vincent Dutordoir, James Hensman, Mark van der Wilk, Carl Henrik Ek, Zoubin Ghahramani, and Nicolas Durand, “Deep neural networks as point estimates for deep gaussian processes,” *NeurIPS*, vol. 34, 2021.
- [10] Jeremiah Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax Weiss, and Balaji Lakshminarayanan, “Simple and principled uncertainty estimation with deterministic deep learning via distance awareness,” *NeurIPS*, vol. 33, pp. 7498–7512, 2020.
- [11] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida, “Spectral normalization for generative adversarial networks,” in *ICLR*, 2018.
- [12] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár, “Focal loss for dense object detection,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
- [13] Ali Rahimi and Benjamin Recht, “Random features for large-scale kernel machines,” *Advances in neural information processing systems*, vol. 20, 2007.
- [14] Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S. Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein, “Deep neural networks as gaussian processes,” in *6th International Conference on Learning Representations*, 2018.
- [15] Joost van Amersfoort, Lewis Smith, Andrew Jesson, Oscar Key, and Yarin Gal, “On feature collapse and deep kernel learning for single forward pass uncertainty,” *arXiv preprint arXiv:2102.11409*, 2021.
- [16] Zhiyun Lu, Eugene Ie, and Fei Sha, “Mean-field approximation to gaussian-softmax integral with application to uncertainty estimation,” *arXiv preprint arXiv:2006.07584*, 2020.
- [17] Jishnu Mukhoti, Viveka Kulharia, Amartya Sanyal, Stuart Golodetz, Philip Torr, and Puneet Dokania, “Calibrating deep neural networks using focal loss,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 15288–15299, 2020.
- [18] Chen Qu, Liu Yang, W Bruce Croft, Johanne R Trippas, Yongfeng Zhang, and Minghui Qiu, “Analyzing and characterizing user intent in information-seeking conversations,” in *The 41st international acm sigir conference on research & development in information retrieval*, 2018, pp. 989–992.
- [19] Gustavo Penha, Alexandru Balan, and Claudia Hauff, “Introducing mantis: a novel multi-domain information seeking dialogues dataset,” *arXiv preprint arXiv:1912.04639*, 2019.
- [20] Ryan Lowe, Nissan Pow, Iulian Serban, and Joelle Pineau, “The ubuntu dialogue corpus: A large dataset for research in unstructured multi-turn dialogue systems,” in *The 16th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, 2015, pp. 285–294.
- [21] Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht, “Obtaining well calibrated probabilities using bayesian binning,” in *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- [22] Agustinus Kristiadi, Matthias Hein, and Philipp Hennig, “Being bayesian, even just a bit, fixes overconfidence in relu networks,” in *International conference on machine learning*. PMLR, 2020, pp. 5436–5446.