# Towards increased trustworthiness of deep learning segmentation methods on cardiac MRI

Jörg Sander[a], Bob D. de Vos[a], Jelmer M. Wolterink[a] and Ivana Išgum[a]

[a]Image Sciences Institute, University Medical Center Utrecht, Utrecht, The Netherlands

## ABSTRACT

Current state-of-the-art deep learning segmentation methods have not yet made a broad entrance into the clinical setting in spite of high demand for such automatic methods. One important reason is the lack of reliability caused by models that fail unnoticed and often locally produce anatomically implausible results that medical experts would not make. This paper presents an automatic image segmentation method based on (Bayesian) dilated convolutional networks (DCNN) that generate segmentation masks and spatial uncertainty maps for the input image at hand. The method was trained and evaluated using segmentation of the left ventricle (LV) cavity, right ventricle (RV) endocardium and myocardium (Myo) at end-diastole (ED) and end-systole (ES) in 100 cardiac 2D MR scans from the MICCAI 2017 Challenge (ACDC). Combining segmentations and uncertainty maps and employing a human-in-the-loop setting, we provide evidence that image areas indicated as highly uncertain, regarding the obtained segmentation, almost entirely cover regions of incorrect segmentations. The fused information can be harnessed to increase segmentation performance. Our results reveal that we can obtain valuable spatial uncertainty maps with low computational effort using DCNNs.

**Keywords:** cardiac MRI segmentation, uncertainty estimation, loss functions, deep learning, Bayesian neural networks

## 1. PURPOSE

Decisions by medical experts are increasingly enriched and augmented by intelligent machines e.g., through computer aided diagnosis (CAD). The quality of the joint decision process would improve if the automatic systems were able to indicate their uncertainty. This assumes that the provided uncertainty information is reliable i.e., valuable to be considered. A system indicating high uncertainty in image areas of incorrect segmentations could be used to detect and subsequently refer these regions to medical experts. Applying such a human-in-the-loop setting would result in increased segmentation performance. In addition, such a setting could mitigate a severe deficiency of current state-of-the-art deep learning segmentation methods which occasionally generate anatomically implausible segmentations[1] that a medical expert would never make. Previous research has mainly focused on the assessment of uncertainty in disease prediction,[2] tissue segmentation[3] and pulmonary nodule detection[4] by utilizing Bayesian neural networks (BNN) or test-time data augmentation techniques.[5] Additional methods to estimate uncertainty are Deep Ensembles[6] and Learned Confidence Estimates.[7] In the former multiple models are trained and the variance of their predictions is used as confidence measure, whereas in the latter the model outputs a confidence measure simultaneously with the prediction.

In this work, using multi-structures segmentation in cardiac MR images, we introduce a method that simultaneously generates segmentation masks and uncertainty maps by using a dilated convolutional network (DCNN). We compare two approaches to obtain uncertainty maps. First, we use entropy maps (e-maps) that can be efficiently generated by any probabilistic classifier as entropy is a theoretically grounded quantification of uncertainty in information theory. Second, we employ Bayesian uncertainty maps (u-maps) that are obtained by Bayesian DCNNs (B-DCNN). In addition, we reveal that a valuable uncertainty measure can be obtained if the applied model is *well calibrated*, i.e. if generated probabilities represent the likelihood of being correct. We demonstrate this by simulating a human-in-the-loop setting and provide evidence that image areas indicated as highly uncertain regarding the obtained segmentation almost entirely cover regions of incorrect segmentations. Hence, the fused information can be employed in clinical practice to inform an expert whether and where the generated segmentation should be adjusted.

---

Send correspondence to J.Sander (email: j.sander@umcutrecht.nl)

## 2. DATA DESCRIPTION

Data from the MICCAI 2017 Challenge on automated cardiac diagnosis (ACDC)[1] was used. The dataset consists of cardiac cine MR images (CMRI) from 150 patients who have been clinically diagnosed in five classes: normal, dilated cardiomyopathy, hypertrophic cardiomyopathy, heart failure with infarction, or right ventricular abnormality. Cases are uniformly distributed over classes. Manual reference segmentations of the left ventricle (LV) cavity, right ventricle (RV) endocardium and myocardium (Myo) at end-diastole (ED) and end-systole (ES) are provided for 100 cases. For each patient, short-axis (SA) CMRIs with 28-40 frames are available, in which the ED and ES frame have been indicated. On average images consist of nine slices where each slice has a spatial resolution of 235×263 voxels (on average). The image slices cover the LV from the base to the apex. In-plane voxel spacing varies from 1.37 to 1.68 mm, with slice thickness from 5 to 10 mm and sometimes inter-slice gap of 5 mm. To correct for differences in voxel size, all 2D image slices were resampled to $1.4 \times 1.4 \, \text{mm}^2$. Furthermore, to correct for intensity differences among images, each MR volume was normalized between [0.0, 1.0] according to the 5th and 95th percentile of intensities in the image. For detailed specifications about the acquisition protocol we refer the reader to Bernard et al.[1]

## 3. METHOD

To perform segmentation of tissue classes in cardiac 2D MR scans, we used the DCNN developed by Wolterink et al.[8] The DCNN architecture comprises a sequence of ten convolutional layers with increasing levels of kernel dilation which results in a receptive field for each voxel of 131×131 voxels, or $18.3 \times 18.3 \, \text{cm}^2$. The network has two input channels which take ED and ES slices as its input. We assume that the network leverages cardiac motion differences between ED and ES time points in order to better localize the target structures. To simultaneously segment the LV, RV, myocardium and background in ED and ES, the network has eight output channels where each channel provides a probability for one of the classes. Softmax probabilities are calculated over the four tissue classes for images acquired in ED and ES. To enhance generalization performance, the model uses batch normalization and weight decay.

To acquire spatial uncertainty maps of the segmentation during testing, two different approaches were evaluated. First, to obtain entropy maps (e-maps) we computed the multi-class entropy per voxel. Second, to obtain Bayesian uncertainty maps (u-maps), we implemented *Monte Carlo dropout* (MC dropout) introduced by Gal & Ghahramani[9] for approximate Bayesian inference. We added dropout as the last operation in all but the final layer (by randomly switching off 10 percent of a layer's hidden units). By enabling dropout during testing, softmax probabilities are obtained with 10 samples per voxel. As an overall measure of uncertainty we used the maximum softmax variance per voxel over all classes. The variance per voxel per class is obtained from the softmax samples for each class. We chose to use the maximum instead of the mean (as e.g., utilized by Leibig et al.[2]) because we found that averaging attenuates the uncertainties.

The quality of e-maps and u-maps depends on the calibration of the acquired probabilities. Previous work[6] revealed that loss functions differ regarding how well the generated probabilities represent the likelihood of being correct. Therefore, we trained the model with three different loss functions: soft-Dice (SD), cross-entropy (CE), and the Brier score (BS),[10] which is equal to the average gap between softmax probabilities and the references. This provides information about accuracy and uncertainty of the model. Computationally the Brier score loss is equal to the squared error between the one-hot encoding of the correct label and its associated probability.

To use four-fold cross-validation we split the dataset into 75 and 25 training and test patients, respectively. Each model is evaluated on the holdout test images and we report combined results for all 100 patients. During training we used images with 151×151 voxel samples, padded to 281×281 to accommodate the 131×131 voxel receptive field. Training samples were augmented by 90 degree rotations of the images and references. The model was trained for 150,000 iterations using the snapshot ensemble technique described in,[11] while after every 10,000th iteration we reset the learning rate to its original value of 0.02 and stored the model. We used mini-batches of size 16 and applied Adam[12] as stochastic gradient descent optimizer. To compare u-maps with e-maps at test time each model was evaluated twice. First, to obtain u-maps we used the last six stored models (iterations 100,000 to 150,000) of each fold to obtain segmentation results. Tissue class per voxel was determined using the mean softmax probabilities over 60 samples (10 samples per voxel per model). In addition, these probabilities

(a) Brier score loss  (b) soft-Dice loss

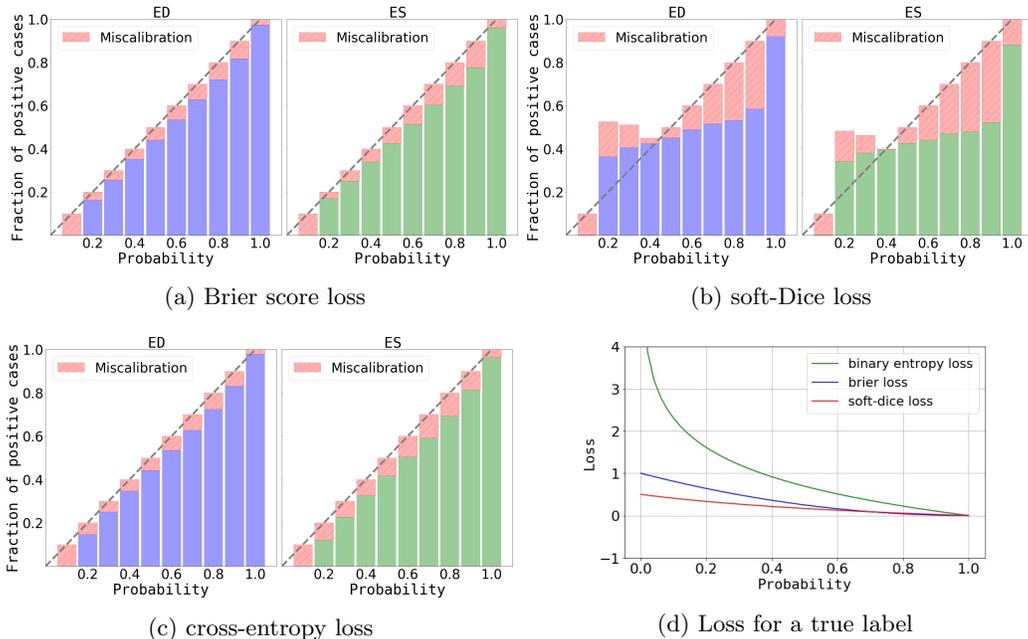(c) cross-entropy loss  (d) Loss for a true label

Figure 1: Reliability diagrams over all tested ED and ES images and tissue classes for Brier, soft-Dice and cross-entropy loss functions. Blue (end-diastole) and green (end-systole) bars quantify the true positive fraction for each probability bin. Red bars quantify the miscalibration of the model where smaller indicates better. If the model is perfectly calibrated, the diagram should match the dashed line.

served to compute the maximum variance (as described in the beginning of this section). Second, to obtain e-maps we solely employed the last stored model of each fold to acquire segmentation results. We disabled dropout during inference and used one forward pass to compute the softmax probabilities and determine the tissue class per voxel. The corresponding e-maps were computed as the entropy in the four-class probability distribution. Finally, for both evaluations as a post-processing step, the 3D probability volumes were filtered by selection of the largest 3D 6-connected component for each class.

The models were implemented using the PyTorch[13] framework and trained on one Nvidia GTX Titan X GPU with 12 GB memory.

## 4. RESULTS AND DISCUSSION

To evaluate whether the obtained per voxel probabilities represent the likelihood of being correct i.e. are well calibrated, we created *Reliability Diagrams*[14] (RD). Figures. 1a, 1b and 1c show the predicted probabilities discretized into ten bins and plotted against the true positive fraction for each bin. If the model is perfectly calibrated, the diagram should match the dashed line. We conclude that a model trained with the soft-Dice loss produces inferior calibrated probabilities compared to the other two loss functions. We conjecture that this could be caused by the relatively low penalty induced by the soft-Dice loss for the model being *underconfident* for true positive tissue labels (see Fig. 1d).

To compare the quality of the obtained uncertainty maps, we simulate a human-in-the-loop setting. We combine the information of predicted segmentation masks with the e-maps or u-maps and assume that voxels above a tolerated uncertainty or entropy threshold are corrected to their reference label by an expert. For each threshold we compute the Dice score for the corrected segmentation mask. Figures 2a and 2b visualize the Dice score as a function of the average percentage of voxels thus referred. We observe a monotonic increase in prediction accuracy when more voxels are referred. E.g., inspecting the referral curves for the Brier score loss in Figure 2b we note that referring on average 1% of the voxels in an image, increases performance for 8, 7 and 5% for RV, Myo and LV, respectively. These results are similar for the u-maps and the e-maps. In each experiment,
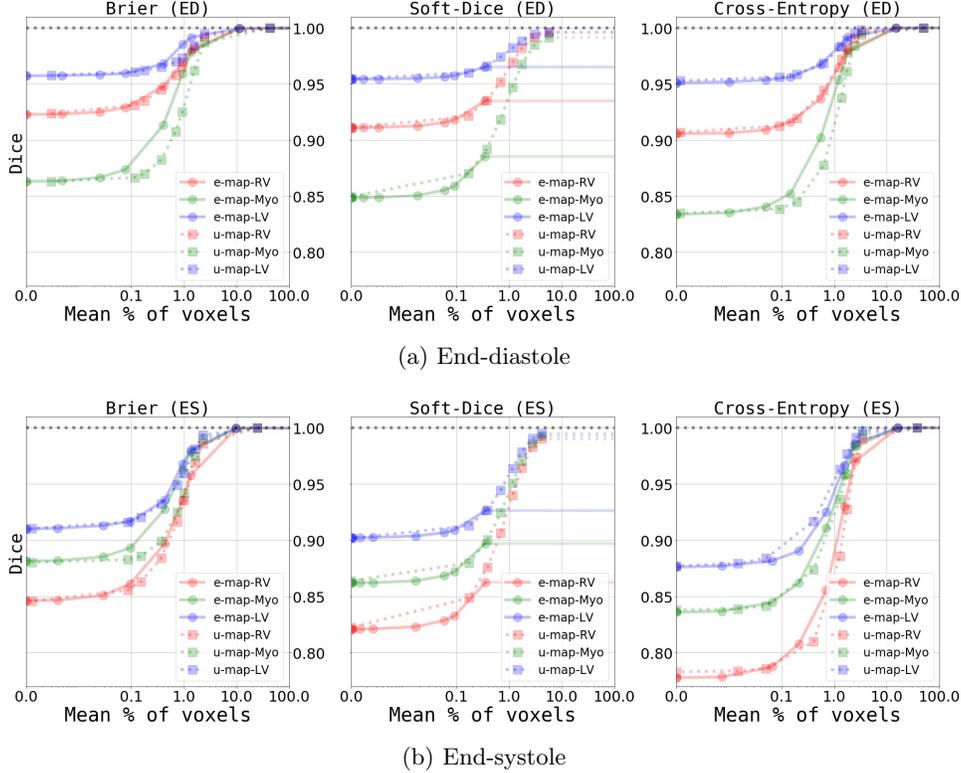
(a) End-diastole



(b) End-systole

Figure 2: Comparison between entropy and Bayesian uncertainty maps for different loss functions (RV in red, myocardium in green and LV in blue). Figures visualize Dice score of the corrected segmentation mask when voxels above a tolerated uncertainty or entropy threshold are corrected to their reference label. x-axis shows mean percentage of voxels referred in an image.

the case in which no voxels are referred for correction is considered the baseline (left most y-axis values). We observe that baseline segmentation performance is highest when the model is trained with the Brier score loss, slightly lower for the soft-Dice, and lowest when cross-entropy is used. Except for the soft-Dice loss we note that u-maps and e-maps follow each other quite closely, which suggests that both carry similar information. Not including the soft-Dice loss, segmentation performance with referral using u-maps or e-maps reaches a Dice score of nearly one when sufficient number of voxels are referred. Hence, we may conclude that areas of uncertainty and entropy almost completely cover the regions of incorrect segmentations[*]. Results obtained after the referral using entropy maps for a model trained with the soft-Dice loss are clearly inferior compared to the performance achieved when using the u-maps. We assume that this is due to the miscalibration of the model (see Figure 1b). Compared to e-maps, u-maps tend to exhibit more uncertainty. This is visually expressed for the cross-entropy loss in Figure 2a, where the Myo referral-curve obtained with u-maps lags behind the corresponding curve that uses the entropy information.

An example result of the segmentation task performed by a model trained with the Brier score loss is shown in Figure 3. The model obviously failed to segment parts of the right ventricle (blue) and we can observe that these errors are covered by entropy and Bayesian uncertainty maps. Figure 4 shows a qualitative comparison of the uncertainty maps for the three different loss functions (corresponding to rows in the figure) that were used during training. Images in the left column visualize the segmentation errors for the three different tissue types using distinct colors. Although we can observe that the performed errors are roughly the same for the different loss functions, we clearly see significant differences between the uncertainty maps. E.g., when inspecting the e-maps (middle column) we notice that errors with respect to the segmentation of the myocardium are not

---

[*]Without covering the complete image in which case all voxels would be referred (corresponding to a trivial solution).

(a) MRI slice      (b) Reference      (c) Automatic      (d) e-map      (e) u-map
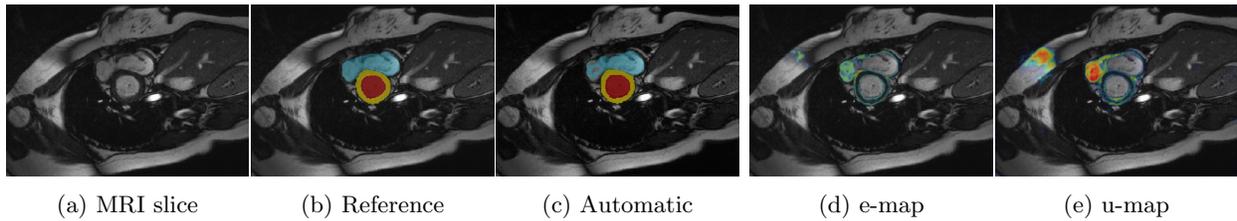
Figure 3: Example of segmentation errors that are covered by high uncertainties. Figure 3a shows the original MRI slice to be segmented. Figure 3b visualizes the manual reference segmentation whereas 3c shows the automatic segmentation mask generated by the model when trained with the Brier score loss. Segmentation errors for the right ventricle (in blue 3b and 3c) are covered by high entropy (figure 3e) and Bayesian uncertainties (figure 3d).



(a) BS: Errors      (b) BS: e-map      (c) BS: u-map

(d) CE: Errors      (e) CE: e-map      (f) CE: u-map

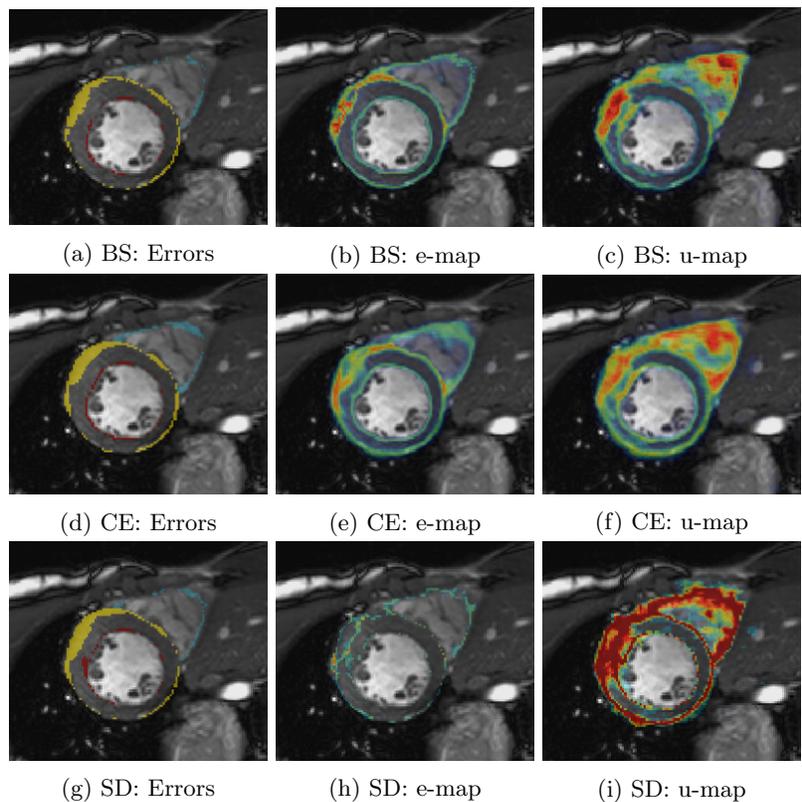(g) SD: Errors      (h) SD: e-map      (i) SD: u-map

Figure 4: Comparison of (left column) segmentation errors of left ventricle (red), myocardium (yellow) and right ventricle (blue); (middle column) Entropy maps; and (right column) Bayesian uncertainty maps for the Brier score (BS), cross-entropy (CE) and soft-Dice (SD) loss (per row). High uncertainties correspond to red and low uncertainties to blue colors.

entirely covered by regions of high uncertainties for a model trained with the soft-Dice loss. In contrast the same regions are almost completely covered by the e-map for a model trained with the Brier score or cross-entropy loss. Furthermore, a model trained with the soft-Dice loss generated u-maps that contain higher uncertainties than u-maps induced by the other two loss functions. We conjecture that this is caused by the miscalibration of the model (see Figure 1b) which has a bias towards generating probabilities that are close to zero or one, leading to large softmax variances per voxels (we used 10 samples per voxel). This does not affect the e-maps because we do not sample predictions for these maps during testing. Besides, the provided examples corroborate our earlier finding that the u-maps contain more uncertain, yet often correctly segmented voxels than the e-maps.

## 5. NEW OR BREAKTHROUGH WORK TO BE PRESENTED

This study shows how automatic segmentation can be combined with spatial uncertainty maps to increase the segmentation performance employing a human-in-the-loop setting. Furthermore, our results reveal that we can obtain valuable spatial uncertainty maps with low computational effort using well-calibrated DCNNs.

## 6. CONCLUSIONS

Using a publicly available cardiac cine MRI dataset, we showed that a (Bayesian) dilated CNN trained with the Brier loss produces valuable Bayesian uncertainty and entropy maps. Our results convey that regions of high uncertainty almost completely cover areas of incorrect segmentations. Well calibrated models enable us to obtain useful spatial entropy maps, which can be used to increase the segmentation performance of the model.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Bernard, O., Lalande, A., Zotti, C., Cervenansky, F., Yang, X., Heng, P.-A., Cetin, I., Lekadir, K., Camara, O., Ballester, M. A. G., et al., "Deep learning techniques for automatic MRI cardiac multi-structures segmentation and diagnosis: Is the problem solved?," *IEEE Transactions on Medical Imaging* (2018).

[2] Leibig, C., Allken, V., Ayhan, M. S., Berens, P., and Wahl, S., "Leveraging uncertainty information from deep neural networks for disease detection," *Scientific reports* **7**(1), 17816 (2017).

[3] Kwon, Y., Won, J.-H., Kim, B. J., and Paik, M. C., "Uncertainty quantification using bayesian neural networks in classification: Application to ischemic stroke lesion segmentation," in [*Medical Imaging with Deep Learning Conference*], (2018).

[4] Ozdemir, O., Woodward, B., and Berlin, A. A., "Propagating uncertainty in multi-stage bayesian convolutional neural networks with application to pulmonary nodule detection," in [*NIPS Workshop on Bayesian Deep Learning*], (2017).

[5] Ayhan, M. S. and Berens, P., "Test-time data augmentation for estimation of heteroscedastic aleatoric uncertainty in deep neural networks," in [*Medical Imaging with Deep Learning Conference*], (2018).

[6] Lakshminarayanan, B., Pritzel, A., and Blundell, C., "Simple and scalable predictive uncertainty estimation using deep ensembles," in [*Advances in Neural Information Processing Systems*], 6402–6413 (2017).

[7] DeVries, T. and Taylor, G. W., "Learning confidence for out-of-distribution detection in neural networks," *arXiv preprint arXiv:1802.04865* (2018).

[8] Wolterink, J. M., Leiner, T., Viergever, M. A., and Išgum, I., "Automatic segmentation and disease classification using cardiac cine MR images," in [*International Workshop on Statistical Atlases and Computational Models of the Heart*], 101–110, Springer (2017).

[9] Gal, Y. and Ghahramani, Z., "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in [*International Conference on Machine Learning (ICML)*], 1050–1059 (2016).

[10] Brier, G. W., "Verification of forecasts expressed in terms of probability," *Monthey Weather Review* **78**(1), 1–3 (1950).

[11] Huang, G., Li, Y., Pleiss, G., Liu, Z., Hopcroft, J. E., and Weinberger, K. Q., "Snapshot ensembles: Train 1, get m for free," *arXiv preprint arXiv:1704.00109* (2017).

[12] Kingma, D. and Ba, J., "Adam: A method for stochastic optimization," in [*ICLR*], **5** (2015).

[13] Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A., "Automatic differentiation in pytorch," (2017).

[14] DeGroot, M. H. and Fienberg, S. E., "The comparison and evaluation of forecasters," *The statistician* , 12–22 (1983).