

1 **Leveraging uncertainty information from deep**
2 **neural networks for disease detection:**
3 **supplementary material**

4 **Christian Leibig^{1,*}, Vaneeda Allken¹, Murat Seçkin Ayhan¹, Philipp Berens^{1,2+}, and**
5 **Siegfried Wahl^{1,3+}**

6 ¹Institute for Ophthalmic Research, Eberhard Karls University, Tübingen, Germany

7 ²Bernstein Center for Computational Neuroscience and Centre for Integrative Neuroscience, Eberhard Karls
8 University, Tübingen, Germany

9 ³Carl Zeiss Vision International GmbH, Germany

10 *christian.leibig@uni-tuebingen.de

11 ⁺Co-senior author

12 **Performance improvement under decision referral for an already trained network**

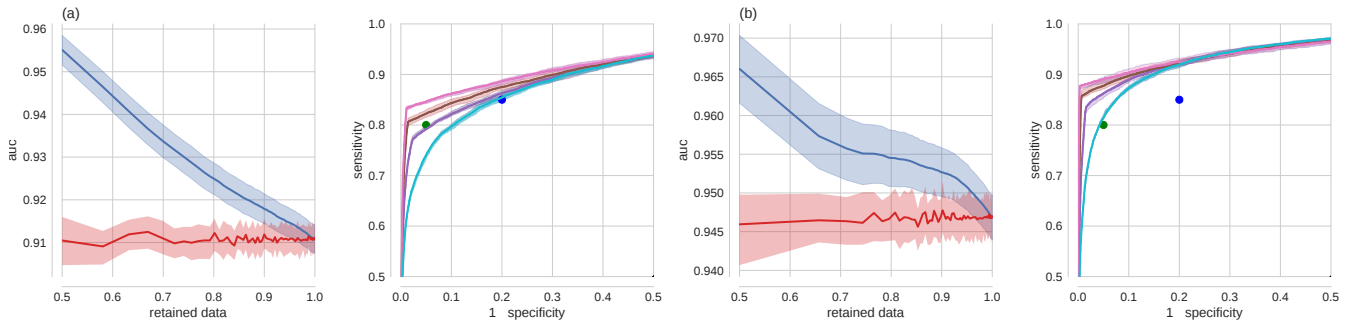


Figure S1. Application of uncertainty-informed decision referral to an existing network (a, left) ROC AUC over the fraction of retained data under uncertainty informed (blue) and random (red) decision referral for the JFnet, recast to detect disease onset 1. (a, right) ROC curves for all data (no referral: turquoise) and different fractions of retained data (90%: purple, 80%: brown, 70%: pink). National UK standards for the detection of sight-threatening diabetic retinopathy (in² defined as moderate DR) from the BDA (80%/95% sensitivity/specificity, green dot) and the NHS (85%/80% sensitivity/specificity, blue dot) are given in all subpanels with ROC curves. (b) same as (a), but for disease onset 2. All subfigures are based on Kaggle DR test images.

13 **Relation between μ_{pred} and σ_{pred}**

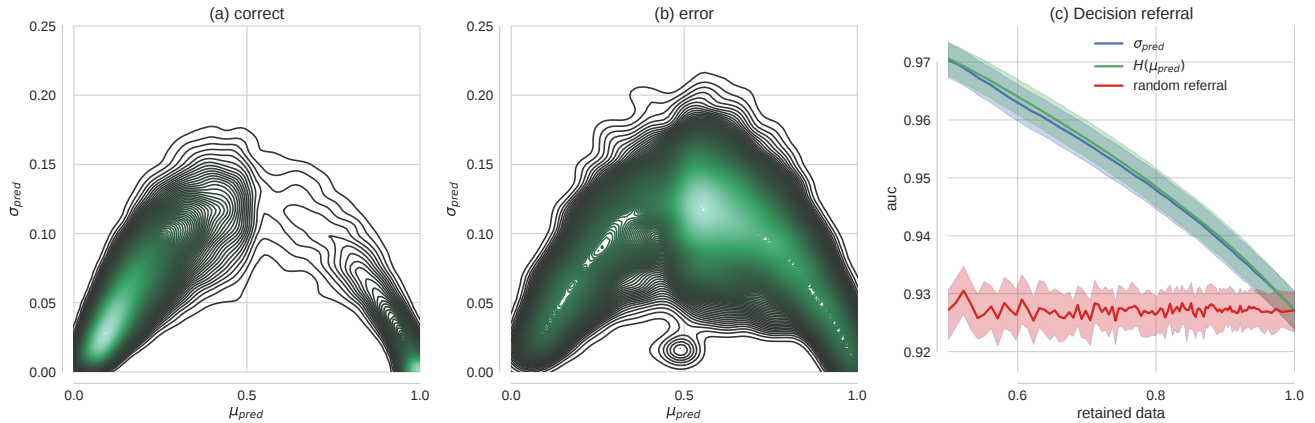


Figure S2. (a) & (b) Relation between first (μ_{pred}) and second (σ_{pred}) moments of approximate predictive posterior for correct (a) and erroneous (b) detection of moderate DR on Kaggle test images. μ_{pred} obtained via MC dropout is more related with σ_{pred} than the network output $p(diseased|image)$ obtained with standard dropout (compare fig. 2). (c) For comparison with our performance results obtained with σ_{pred} as uncertainty (fig. 4, S1), we quantified the uncertainty in terms of the binary entropy $H(p) = -(p \log p + (1 - p) \log (1 - p))$ as an alternative uncertainty measure which is applicable to both the Bayesian and conventional network output. When using $H(\mu_{pred})$ (green curve) we were able to achieve similar performance improvements under decision referral compared to when using σ_{pred} (blue curve) instead. Random referral is shown in red.

14 **Gaussian processes**

15 Gaussian processes (GPs) enable probabilistic kernel machines for solving regression and classification problems. The GP
 16 inference takes place in a function space and a kernel is a covariance function $k(\mathbf{x}^i, \mathbf{x}^j)$ that estimates the covariance of two
 17 latent variables $f(\mathbf{x}^i)$ and $f(\mathbf{x}^j)$ in terms of input vectors³. Given a *GP-prior*, e.g., $\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$, over a set of latent variables,
 18 where \mathbf{K} is a covariance matrix determined by the choice of covariance function $k(\cdot, \cdot)$, a *GP-posterior* $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$, where \mathbf{y} is the
 19 target vector and \mathbf{X} is the data matrix, can be inferred in the light of the observed data, namely the likelihood $p(\mathbf{y}|\mathbf{f})$. In the
 20 process, the hyperparameters of the covariance function are also automatically tuned via likelihood maximization³.

21 GP classification is essentially *binary* and class labels are in $\{-1, +1\}$. Similar to the case of the *logistic regression*, the
 22 goal is to assign $p(y^* = +1|\mathbf{x}^*)$ for an unseen test example \mathbf{x}^* . But, the underlying Bayesian treatment of GPs results in an
 23 *averaged predictive probability*³, instead of a single point estimate, for \mathbf{x}^* via Eqs.1 and 2.

$$\bar{\lambda}(f^*) = \int \lambda(f^*)p(f^*|\mathbf{X}, \mathbf{y}, \mathbf{x}^*)df^*, \text{ where} \quad (1)$$

24

$$p(f^*|\mathbf{X}, \mathbf{y}, \mathbf{x}^*) = \int p(f^*|\mathbf{X}, \mathbf{x}^*, \mathbf{f})p(\mathbf{f}|\mathbf{X}, \mathbf{y})d\mathbf{f}, \quad (2)$$

25 $\lambda(f^*)$ is a squashing function that maps its inputs into $[0, 1]$. A common choice of $\lambda(f^*)$ is in Eq. 3.

$$p(y^* = +1|\mathbf{x}^*) = \lambda(f^*) = \frac{1}{1 + \exp(-f^*)}. \quad (3)$$

26 The exact computation of the averaged predictive probability in Eq.2 is intractable and approximation methods are used
 27 to this end. A comprehensive review⁴ of approximate inference methods for GP classification is available in literature. Two
 28 commonly used methods are Laplace Approximation (LA)³ and Expectation Propagation (EP)⁶.

29 LA replaces $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ with a Gaussian approximation centered at the mode of $p(\mathbf{f}|\mathbf{X}, \mathbf{y})$ ³. LA is a local method in the sense
 30 that it exploits the properties of the posterior at a particular location only, e.g., at the posterior distribution's mode⁴. Due to its
 31 simplicity, LA is very fast; however, it may lead to substantial underestimations of the mean and covariance, especially, in
 32 high-dimensional spaces because the mode and mean can be far from each other^{4,5}.

33 The EP is a general approximation tool that can perform approximate inference more accurately than other methods like
 34 Monte Carlo, LA and variational Bayes⁶. Thus, it is heavily used for GP learning. EP is a global method in the sense that it
 35 utilizes many local approximations in order to approximate the posterior^{3,4}. The use of many local approximations results
 36 in a small global divergence between the posterior and its proxy⁴. As a result, EP delivers accurate marginals, reliable class
 37 probabilities and faithful model selection⁴. On the downside, the convergence of EP is not generally guaranteed and its runtime
 38 is 10 times longer in comparison with LA⁴.

39 In summary, a GP classification can be made by just checking the sign of the predictive function sampled from Eq. 2. In this
 40 regard, if one only cares about the error rate or the computational resources are limited, LA may be a practical solution⁴. On
 41 the other hand, EP offers higher quality approximations and it should be considered first, despite the associated computational
 42 costs and risk of not converging, when the quality of approximation matters. In the face of practical challenges, such as
 43 non-convergence or exceptionally long runtimes, one can always fall back onto LA.

44 **Neural Network Covariance Function**

45 A shallow neural network with infinitely many hidden units and Gaussian weights, such as $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma)$, converges into a GP
 46 that can be constructed with the following covariance function^{3,7,8}:

$$k(\mathbf{x}, \mathbf{z}) = \frac{2}{\pi} \sin^{-1} \left(\frac{2\tilde{\mathbf{x}}^T \Sigma \tilde{\mathbf{z}}}{\sqrt{(1 + 2\tilde{\mathbf{x}}^T \Sigma \tilde{\mathbf{z}})(1 + 2\tilde{\mathbf{z}}^T \Sigma \tilde{\mathbf{x}})}} \right), \quad (4)$$

47 where $\tilde{\mathbf{x}} = (1, \mathbf{x})^T$ is an augmented input vector and Σ is a weight covariance matrix. In other words, a GP with the neural
 48 network covariance function (Eq. (4)) emulates a shallow network.

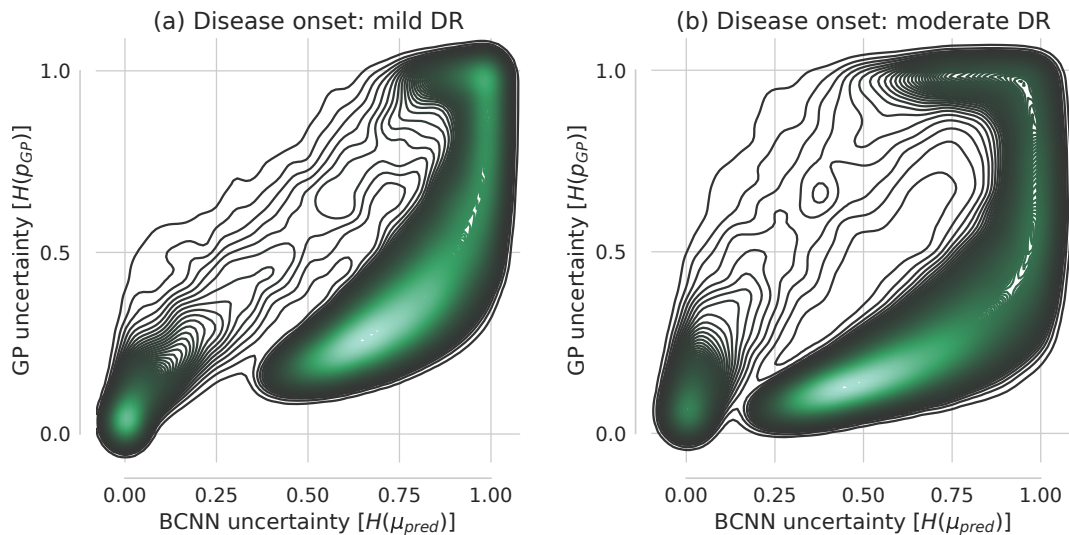


Figure S3. Comparison of Bayesian network (MC dropout) vs. Gaussian process (GP) uncertainty. (a) & (b) Density plots for uncertainties [entropy] obtained from BCNNs vs. GPs for detecting mild (a) and (b) moderate DR from Kaggle data respectively. For the majority of the data, the BCNN uncertainties tend to be larger than the uncertainties from the GP predictions.

References

- 49 **1.** Gal, Y. *Uncertainty in Deep Learning*. Ph.D. thesis, University of Cambridge (2016).
- 50 **2.** Younis, N., Broadbent, D. M., Harding, S. P. & Vora, J. P. Incidence of sight-threatening retinopathy in Type 1 diabetes in
- 51 a systematic screening programme. *Diabetic medicine : a journal of the British Diabetic Association* **20**, 758–765 (2003).
- 52 **3.** Rasmussen, Carl Edward, *Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning)* 2005,
- 53 The MIT Press, 026218253X.
- 54 **4.** Nickisch, Hannes, Rasmussen, Carl Edward, Approximations for binary Gaussian process classification, *Journal of*
- 55 *Machine Learning Research*, **9**, 2035–2078, (2008).
- 56 **5.** Kuss, Malte, Rasmussen, Carl Edward, Assessing approximate inference for binary Gaussian process classification,
- 57 *Journal of Machine Learning Research*, **6**, 1679–1704, (2005).
- 58 **6.** Minka, Thomas P. *A family of algorithms for approximate Bayesian inference*, Ph.D. thesis, Massachusetts Institute of
- 59 Technology (2001).
- 60 **7.** Neal, Radford M., *Bayesian Learning for Neural Networks*, Springer-Verlag New York, Inc., Secaucus, NJ, USA,
- 61 0387947248, (1996)
- 62 **8.** Williams, Christopher K. I. and Barber, David, Bayesian Classification With Gaussian Processes, *IEEE Trans. Pattern*
- 63 *Anal. Mach. Intell.*, **20**, 1342–1351, (1998).
- 64