

PARAMETRIC DENSITY ESTIMATION WITH UNCERTAINTY USING DEEP ENSEMBLES

Anonymous authors

Paper under double-blind review

ABSTRACT

1 In parametric density estimation, the parameters of a known probability density
2 are typically recovered from measurements by maximizing the log-likelihood.
3 Prior knowledge of measurement uncertainties is not included in this method, po-
4 tentially producing degraded or even biased parameter estimates. We propose
5 an efficient two-step, general-purpose approach for parametric density estimation
6 using deep ensembles. Feature predictions and their uncertainties are returned
7 by a deep ensemble and then combined in an importance weighted maximum
8 likelihood estimation to recover parameters representing a known density along
9 with their respective errors. To compare the bias-variance tradeoff of different
10 approaches, we define an appropriate figure of merit. We illustrate a number of
11 use cases for our method in the physical sciences and demonstrate state-of-the-art
12 results for X-ray polarimetry that outperform current classical and deep learning
13 methods.

14 1 INTRODUCTION

15 The majority of state-of-the-art NN performances are single (high-dimensional) input, multiple-
16 output tasks, for instance classifying images (Krizhevsky et al., 2012), scene understanding (Red-
17 mon et al., 2015) and voice recognition (Graves et al., 2006). These tasks typically involve one input
18 vector or image and a single output vector of predictions.

19 In parametric density estimation, there is a known probability density that the data (or latent features
20 of the data) are expected to follow. The goal is to find representative distribution parameters for a
21 given dataset. In simple cases where the likelihood is calculable, maximum likelihood estimation
22 can be used effectively. In cases where latent features of the data follow a known distribution (e.g.,
23 heights of people in a dataset of photographs), NNs can potentially be used to directly estimate the
24 distribution parameters. For clarity, we define this direct/end-to-end approach as parametric feature
25 density estimation (PFDE). Such an approach requires employing entire datasets (with potentially
26 thousands to millions of high-dimensional examples) as inputs in order to output a vector of den-
27 sity parameters. Furthermore, to be useful these NNs would need to generalize to arbitrarily sized
28 dataset-inputs.

29 One example of NNs making sense of large dataset-inputs is found in natural language processing.
30 Here large text corpora, converted to word vectors (Pennington et al., 2014; Devlin et al., 2019),
31 can be input and summarized by single output vectors using recurrent neural networks (RNNs), for
32 instance in sentiment analysis (Can et al., 2018). However, these problems and RNNs themselves
33 contain inductive bias – there is inherent structure in text. Not all information need be given at once
34 and a concept of memory or attention is sufficient (Vaswani et al., 2017). The same can be said
35 about time domain problems, such as audio processing or voice recognition. Memory is inherently
36 imperfect – for PFDE, one ideally wants to know all elements of the ensemble at once to make
37 the best prediction: sequential inductive bias is undesirable. Ultimately, memory and architectural
38 constraints make training NNs for direct PFDE computationally intractable.

39 On the other hand, density estimation on data directly (not on its latent features), is computationally
40 tractable. Density estimation lets us find a complete statistical model of the data generating process.
41 Applying deep learning to density estimation has advanced the field significantly (Papamakarios,
42 2019). Most of the work so far focuses on density estimation where the density is unknown *a priori*.
43 This can be achieved with non-parametric methods such as neural density estimation (Papamakarios

44 et al., 2018), or with parametric methods such as mixture density networks (Bishop, 1994). In PFDE,
 45 however, we have a known probability density over some features of the whole dataset. The features
 46 may be more difficult to predict accurately in some datapoints than others.

47 Typical parametric density estimation does not make use of data uncertainties where some elements
 48 in the dataset may be more noisy than others. Not including uncertainty information can lead to
 49 biased or even degraded parameter estimates. The simplest example of parametric density estimation
 50 using uncertainties is a weighted mean. This is the result of a maximum likelihood estimate for a
 51 multi-dimensional Gaussian. For density estimation on predicted data features, PFDE, we would
 52 like a way to quantify the predictive uncertainty. A general solution is offered by deep ensembles
 53 (Lakshminarayanan et al., 2017). While these are not strictly equivalent to a Bayesian approach,
 54 although they can be made such using appropriate regularization (Pearce et al., 2018), they offer
 55 practical predictive uncertainties, and have been shown to generalize readily (Fort et al., 2019).
 56 [Additionally Ovadia et al. \(2019\) have shown deep ensembles perform the best across a number](#)
 57 [of uncertainty metrics, including dataset shift, compared to competing methods such as stochastic](#)
 58 [variational inference and Monte Carlo methods.](#)

59 In this work, we propose a NN approach that circumvents large dataset-input training or recurrent
 60 architectures to predict known feature density parameters over large input datasets. We use predic-
 61 tive uncertainties on features of individual dataset elements as importance weights in a maximum
 62 likelihood estimation. We will show that estimating known density parameters in a 2-step approach
 63 provides greater interpretability and flexibility. We are able to predict uncertainties on our density
 64 parameter estimates using bootstrap methods (Efron, 1979). Our method is widely applicable to a
 65 number of applied machine learning fields; §3 showcases a few important examples.

66 **Contributions:** Our contributions in this paper are as follows: (1) We introduce a general, flexi-
 67 ble method for PFDE using NNs. The method can be applied to any domain requiring PFDE. We
 68 illustrate a number of varied domain examples in the physical sciences in §3. (2) In an in-depth
 69 evaluation we show that our method outperforms not only classical methods for density estimation,
 70 but also standard NN implementations in an application to X-ray polarimetry. (3) We investigate the
 71 bias-variance tradeoff associated with our method and introduce a tuneable hyperparameter to con-
 72 trol it. *Note:* In the following we focus on regression examples, (since unbinned density estimation
 73 is preferable to binned). However, a similar method can be applied to prediction examples where
 74 softmax class probabilities are used as heteroscedastic aleatoric uncertainty.

75 2 IMPORTANCE WEIGHTED ESTIMATION WITH DEEP ENSEMBLES

76 2.1 PROBLEM SETUP AND HIGH-LEVEL SUMMARY

77 We wish to estimate the feature density parameters of N high dimensional data points $\{\mathbf{x}\}$:
 78 $f(\{\mathbf{x}_n\}_{n=1}^N)$. Here $\mathbf{x} \in \mathbb{R}^D$ can be any high dimensional data (e.g. images, time series). N is
 79 [arbitrary, although usually large since otherwise density estimation is inaccurate](#). For example, con-
 80 sider estimating the mean and variance of human heights from a dataset consisting of photographs
 81 of people. A person’s height in each photograph is the image feature and we know this feature
 82 approximately follows a Gaussian distribution. [We develop a method that can estimate the density](#)
 83 [parameters \(mean and variance\) and generalize to any dataset of photographs.](#)

84 In general, the function f mapping the high dimensional data points to the desired density paramet-
 85 ers is unknown, since the high dimensional data is abstracted from its features. [Learning \$f\$ directly](#)
 86 [is typically infeasible because an entire ensemble of inputs \$\{\mathbf{x}_n\}_{n=1}^N\$ must be processed simultane-](#)
 87 [ously to estimate density parameters, and this approach would have to generalize to arbitrary \$N\$ and](#)
 88 [density parameter values. We discuss some special cases where this is possible in §1.](#) However, the
 89 function g mapping data features y_n to the density parameters is known.

90 We cast this as a supervised learning problem where we have a dataset D consisting of N data points
 91 $D = \{\mathbf{x}_n, y_n\}_{n=1}^{N_{\text{train}}}$ with labels $y \in \mathbb{R}^K$ where $\mathbf{x} \in \mathbb{R}^D$. We want to estimate the density parameters
 92 $\psi_1, \psi_2, \dots, \psi_k$ for an unseen test set $g(\{y_n\}_{n=1}^{N_{\text{test}}})$ for arbitrary N_{test} .

93 The basic recipe that comes to mind is training a single NN to predict output labels $\{y_n\}_{n=1}^N$ then
 94 evaluate g directly. This ignores the high variance in single NN predictions (dependent on train-

95 ing/random initialization), that some individual examples may be more informative than others, and
 96 that an objective to predict the most accurate output labels may not be the best for predicting good
 97 density parameters (high bias may be introduced, for instance).

98 Our hybrid approach is as follows. (i) Train a deep ensemble of M NNs¹ to predict $\{y_n, \sigma_n\}_{n=1}^N$
 99 where σ_n is the total uncertainty on each prediction y_n , (ii) use the $\{\sigma_n\}_{n=1}^N$ as weights in an
 100 importance weighted maximum likelihood estimate. The next section, §2.2, describes procedure (i).

101 2.2 DEEP ENSEMBLES

102 Deep ensembles (Lakshminarayanan et al., 2017) return robust and accurate supervised learning pre-
 103 dictions and predictive uncertainties, which enable the best density parameter predictions. These use
 104 an ensemble of individual NNs (with different random initializations) trained to predict features and
 105 their aleatoric uncertainties. Final predictions and their epistemic uncertainties are then recovered
 106 by combining the estimates from each of the NNs in the ensemble.

107 In regression, deep ensembles model heteroscedastic aleatoric σ_a uncertainty by modifying the typi-
 108 cal mean-squared errors (MSE) objective to a negative log-likelihood (NLL) (Lakshminarayanan
 109 et al., 2017),

$$\text{Loss}(y|\mathbf{x}) = \frac{1}{2} \log \sigma_a^2(\mathbf{x}) + \frac{1}{2\sigma_a^2(\mathbf{x})} \|y - \hat{y}(\mathbf{x})\|_2^2. \quad (1)$$

110 Extensions using more complex distributions like mixture density networks or heavy tailed distribu-
 111 tions may be more applicable to certain problems with prior knowledge about the error distribution.
 112 In practice, the log-likelihood of any exponential family could be used; we find this simple Gaussian
 113 approach to be sufficient and robust for regression problems. Our results in §3.4 for a compare a
 114 Gaussian and Von Mises distribution.

115 Epistemic uncertainty σ_e is modelled using a uniformly weighted ensemble of M NNs each
 116 trained starting from a different random initialization. The regression prediction and uncertainty
 117 are approximated by the mean and standard deviation over the M NN ensemble predictions re-
 118 spectively (each NN in the ensemble contributes equally) i.e. $\hat{y}(\mathbf{x}) = M^{-1} \sum_{m=1}^M \hat{y}_m(\mathbf{x})$ and
 119 $\sigma_e^2(\mathbf{x}) = \text{Var}(\{\hat{y}_m(\mathbf{x})\}_{m=1}^M)$. The epistemic uncertainty is then combined with the aleatoric in
 120 quadrature to arrive at the total uncertainty: $\sigma^2 = \sigma_a^2 + \sigma_e^2$. Typically $M \sim 5 - 15$.

121 In part (i) of our hybrid approach for PFDE, we train a deep ensemble to minimize the NLL (1) on
 122 desired features y . We follow the deep ensemble training procedure outlined in Lakshminarayanan
 123 et al. (2017) (with recast loss function from Kendall & Gal (2017)) without using adversarial examples,
 124 using the full dataset for each NN. Since the individual density parameters over predicted features
 125 are the final desired values in PFDE, it is possible that an objective maximizing feature accuracy on
 126 the validation set is not the true objective. This is possible if the training dataset is biased or the
 127 model (1) is highly misspecified for the particular problem. The Kitaguchi et al. (2019) single CNN
 128 method in table 1, §3.4, shows a clear case of training bias. If de-biasing the training dataset or using
 129 a more appropriate model is not possible, we have identified two potential ways of ameliorating this
 130 issue for PFDE:

- 131 1. Include terms in the individual NN objectives to penalize known sources of bias.
- 132 2. Select the top M performing NNs, as measured by a criterion that includes density param-
 133 eter prediction bias on a held out test set.

134 In practice both can be used simultaneously. However, the former runs into batch size problems
 135 (since one needs a large sample size to accurately estimate bias), and the source of bias is not always
 136 well understood. The latter naturally arises from the use of deep ensembles, but could include its
 137 own unwanted bias and risk underestimating the epistemic uncertainty. We compare selecting the
 138 top performing NNs for the ensemble by a domain specific criterion against randomly selecting NNs
 139 for the ensemble in §3.

¹We note that the NN architecture used will of course depend on the dataset domain.

140 2.3 IMPORTANCE WEIGHTED LOG-LIKELIHOOD

141 Provided a mapping between high dimensional inputs and interpretable features $\mathbf{x}_n \mapsto y_n$, we can
 142 calculate the density parameters $\psi_1, \psi_2, \dots, \psi_k$ by minimizing the appropriate negative log-likelihood
 143 function $p(\{y_n\}|\psi_1, \psi_2, \dots, \psi_k)$. Some feature predictions y_n will have greater total predictive uncer-
 144 tainties, σ_n . We estimate feature density parameters by incorporating the total uncertainty into an
 145 importance weighted maximum likelihood estimate. This makes up part (ii) of our hybrid method.

146 An importance weight quantifies the relative importance of one example over another. Importance
 147 weighting an element should be the same as if that element were included multiple times in the
 148 dataset, proportional to its importance weight Karampatziakis & Langford (2011). The deep ensemble,
 149 once trained, will act as mapping between high dimensional inputs \mathbf{x}_n and feature-uncertainty
 150 output pairs y_n, σ_n . For each input \mathbf{x}_n there will be M output pairs $\{\hat{y}_{nm}, (\sigma_a)_{nm}\}_{m=1}^M$, one for
 151 each NN in the deep ensemble. Both the features \hat{y}_{nm} and aleatoric uncertainty variances $(\sigma_a)_{nm}^2$
 152 can be combined by taking the appropriate mean over m ; this mean may depend on the distribution
 153 used in (1), but for the simple Gaussian case the standard mean is sufficient. Taking the mean results
 154 in a single output pair $(\hat{y}_n, (\sigma_a)_n)$ for each input. Epistemic uncertainties are included as in §2.2,
 155 resulting in the final output (\hat{y}_n, σ_n) .

156 In order to use all possible information when estimating the desired density parameters $\psi_1, \psi_2, \dots, \psi_k$,
 157 we define an importance weighted negative log-likelihood function

$$L_w(\{\hat{y}_n\}, \psi_1, \psi_2, \dots, \psi_k) = - \sum_{n=1}^N w_n \log \mathcal{L}(\hat{y}_n | \psi_1, \psi_2, \dots, \psi_k), \quad (2)$$

158

$$w_n = \sigma_n^{-\lambda} \quad (3)$$

159 Each individual prediction y_n has an associated importance weight w_n . The $\sigma_n^{-\lambda}$ term weights each
 160 y_n by its predictive uncertainty. The hyperparameter $\lambda \geq 0$ controls the importance weighting distri-
 161 bution. A high λ means the y_n with the lowest (estimated) MSE will dominate the final ensemble
 162 statistic. As always in estimation problems, there is a trade-off between lower variance predictions
 163 and more bias. This can be tuned for a specific application using λ ; we discuss the procedure in
 164 detail in our example application, §3. Final density parameters are found by minimizing (2) over the
 165 domain of the density parameters ψ .

166 Typically, the weights in weighted likelihood estimation are determined heuristically (Hu & Zidek,
 167 2002). In this example, we choose $w = \sigma^{-\lambda}$ since it approximates the simple functional form of
 168 the likelihood used in a weighted mean estimate ($\lambda = 2$). This weighting choice is also inspired
 169 by the dispersion parameter used in generalized linear models (GLMs) (Nelder & Wedderburn,
 170 1972). We expect that this weighting will retain similar robustness properties in terms of model
 171 fitting, and will generalize well to many domains. However, of course, any decreasing function
 172 $f : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ may be used to determine weights, with the most suitable choice of function f
 173 within a given class of functions (in our case, parameterized by λ) to be determined by either cross-
 174 validation or performance on a holdout set. In some applications it is possible to find the exact
 175 weighting function [in prep., reference deleted to maintain integrity of review process]. Further
 176 discussion of weight choice in our application is given in section §3.4.

177 Confidence intervals on the density parameters can be calculated using the non-parametric bootstrap
 178 Efron (1979): select N y_n, σ_n pairs with replacement and minimize (2). In the limit of many trials
 179 with different random subsamples, this will give the output distribution on the density parameters.

180 2.4 DENSITY PARAMETER REGRESSION

181 For a special class of parameterized densities it is possible to find the global minimizer or minimize
 182 (2) analytically (e.g. for a multivariate Gaussian). In practice, the majority of parametric densities
 183 of interest for PFDE are likely to be convex (exponential families, our application example §3, etc.),
 184 so will fall into this special class. In the general case, minimization is performed numerically to find
 185 locally optimal solutions.

186 In this work, we employ Ipopt (Wächter & Biegler, 2006), an open-source interior-point solver for
 187 large-scale non-convex optimization problems, to minimize (2). This method can be used for convex

188 or non-convex parametric density estimates, but only convex ones are guaranteed to be global opti-
 189 mal. Because Ipopt finds locally optimal solutions, which are highly dependent upon an initial guess
 190 of the parameters provided to the solver, in the non-convex case, we recommend nested sampling
 191 Feroz et al. (2009) to test many initial guesses and then select the best local solution. Constraints
 192 on the density parameters, for instance if they have a finite domain, can be incorporated for both the
 193 convex and non-convex case. **Of course, any optimizer appropriate for (2) can be used and this will**
 194 **depend on the problem.**

195 The overall training and evaluation procedure is summarized in Algorithm 1.

Algorithm 1: Pseudocode for our PFDE method.

- 1: Identify output features y_n relevant to the desired density parameter(s) (e.g., subject height in photographs).
 - 196 2: Train a deep ensemble of NNs using loss function (1) to maximise accuracy on the desired output features
 - 3: Evaluate the density parameter(s) using importance weights by minimizing (2).
 - 4: Tune λ hyperparameter for the specific application.
-

197 3 EXPERIMENTS

198 3.1 X-RAY POLARIMETRY

199 Measuring X-ray polarization has been a major goal in astrophysics for the last 40 years. X-ray po-
 200 larization can provide essential measurements of magnetic fields very close to high energy sources,
 201 such as accreting black holes and astrophysical jets (Weisskopf, 2018). The recent development
 202 of photoelectron tracking detectors (Bellazzini et al., 2003) has greatly improved the prospects of
 203 doing so. X-ray polarization telescopes with photoelectron tracking detectors directly image elec-
 204 tron tracks formed from photoelectrons scattered by the incoming X-ray photons. We describe an
 205 application of our hybrid PFDE method to X-ray polarimetry using photoelectron tracking detec-
 206 tors. We use data from the upcoming NASA Imaging X-ray Polarization explorer (IXPE) (Sgrò &
 207 IXPE Team, 2019) as a working example. The problem of recovering polarization parameters from
 208 a dataset of (IXPE) electron track images has recently been announced as an open problem in the
 209 machine learning community (Moriakov et al., 2020).

210 The linear polarization of light can be fully described by two degrees of freedom: the polarization
 211 fraction $0 \leq \Pi \leq 1$, (0% – 100%), and the electric vector position angle $-\pi/2 \leq \phi \leq \pi/2$.
 212 These can be thought of as the magnitude and direction of a vector perpendicular to the direction
 213 of propagation of the light. In imaging X-ray polarimetry, when the detector images an X-ray
 214 source, it measures individual 2D images of electron tracks excited by incoming X-ray photons.
 215 The initial directions the electrons travel follow a known probability density that depend on the
 216 source polarization, and the problem is to recover the polarization parameters Π and ϕ from the
 217 collected dataset of 2D track images.

218 In the case of IXPE, charge tracks are imaged by hexagonal pixels. Fig. 1 shows some example
 219 photoelectron tracks at different X-ray energies. Each track represents the interaction of a single
 220 photon with a single gas molecule. The initial track angle y follows the probability density

$$p(y | \Pi, \phi) = \frac{1}{2\pi} (1 + \Pi \cos(2(y + \phi))) , \quad (4)$$

221 where Π and ϕ are fixed polarization parameters that depend on the source. By estimating y for a
 222 large number of tracks, we may recover the original polarization parameters Π and ϕ , using para-
 223 metric density estimation.

224 Track morphologies vary greatly with energy (and even for the same energy); this affects how dif-
 225 ficult it is to recover an accurate initial photoelectron angle y . Low energy tracks are typically less
 226 elliptical and so more difficult to estimate. For this reason it is essential to incorporate some form of
 227 quality control in the tracks used for polarization estimates.

228 Current IXPE methods estimate individual track y using a moment analysis (Sgro, 2017). This
 229 calculates the first, second and third charge moments using the 2D coordinates of the hexagonal

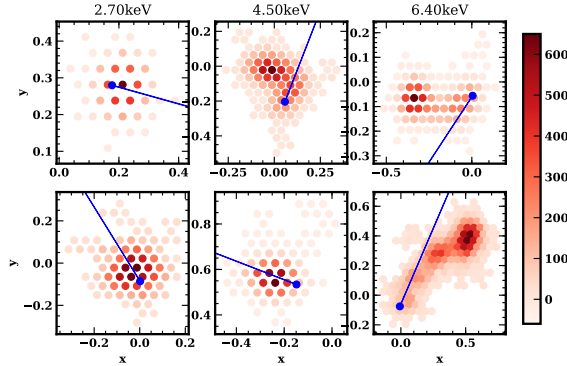


Figure 1: Example IXPE track images at 2.7, 4.5 and 6.4 keV energies (columns). The blue lines show the initial photoelectron direction; the angle of these lines is y . Color represents the amount of charge deposited in a hexagonal pixel. Track morphology (and thus angle reconstruction) depends strongly on energy.

230 detector pixels, combining them to extract y . For each track, a single $-\pi \leq y \leq \pi$ is output.
 231 The polarization parameters are then estimated using a standard (unweighted) MLE. The moment
 232 analysis additionally outputs an estimate of the track ellipticity, which can be used as a proxy for
 233 y estimation accuracy. The standard moment analysis uses a track cut to improve polarization re-
 234 recovery – 20% of the tracks are cut based on ellipticity. NNs have also recently been applied to
 235 this problem Kitaguchi et al. (2019). This approach uses single CNNs for classification on y , with
 236 binned fits to y histograms to extract polarization parameters and track quality cuts. Our hybrid
 237 method exhibits significantly improved performance over both the standard IXPE method and this
 238 basic NN approach.

239 3.2 PARAMETRIC FEATURE DENSITY ESTIMATION

240 Following §2, we define CNNs that take single track images as input and $(\hat{y}, \hat{\sigma})$ as output. **In this**
 241 **case the track angles y are the data features that follow the known density (4), the density parameters**
 242 **$\Pi \equiv \psi_1, \phi \equiv \psi_2$, and the CNNs will make up the deep ensemble.**

243 To make the hexagonal track images admissible inputs to standard CNN architectures, we first
 244 convert the hexagonal images to square image arrays by shifting every other column and rescaling
 245 the distance between points, as described in Steppa & Holch (2019). Since there are two possible
 246 shifts (odd and even rows), we apply both and stack the two shifted images, similar to color channels
 247 in rgb images. We do this to more closely approximate spatial equivariance of the CNN convolution
 248 kernels in the hexagonal space. At test time, we apply the deep ensemble to the same track 3 times,
 249 each time rotated by 120° in hexagonal space. We find this reduces all relevant prediction bias on \hat{y}
 250 (and later Π, ϕ) introduced when converting from hexagonal to square coordinates.

251 To recover Π, ϕ we need to predict $2y$, so we use the loss function (1) but parameterize the true angle
 252 y as a 2D vector $\mathbf{v} = (\cos 2y, \sin 2y)$ to capture the periodicity. The loss function is as follows:

$$\text{Loss}(\mathbf{v}, \hat{\mathbf{v}}) = \frac{1}{2} \log \hat{\sigma}^2 + \frac{1}{2\hat{\sigma}^2} \|\mathbf{v} - \hat{\mathbf{v}}\|_2^2. \tag{5}$$

253 The final NN ensembles output the 3-vector $(\hat{\mathbf{v}}, \hat{\sigma})$. **In this case the mean over ensemble predictions**
 254 **is calculated using the circular mean of $\{\hat{\mathbf{v}}_m\}_{m=1}^M$.** Then $\hat{y} = \frac{1}{2} \arctan \frac{\hat{v}_2}{\hat{v}_1}$. To calculate the final
 255 Π, ϕ with an ensemble of M NNs for a given test dataset with N tracks we minimize the importance
 256 weighted NLL (2) with likelihood

$$\mathcal{L}(\hat{y}_n | \Pi, \phi) = \frac{1}{2\pi} (1 + \Pi \cos(2(\hat{y}_n + \phi))). \tag{6}$$

257 We can recast this as the convex optimization problem

$$\begin{aligned} \underset{\mathbf{x}}{\text{minimize}} \quad & - \sum_{n=1}^N \hat{\sigma}_n^{-\lambda} \log(1 + \mathbf{v}_n^T \mathbf{x}) \\ \text{subject to} \quad & \|\mathbf{x}\|_2 \leq 1 \end{aligned} \tag{7}$$

Energy	Moments	Mom. w/ cuts	Kitaguchi et al.	Single	Ensemble	IW Ensemble (Random)	IW Ensemble (Top MSE)	IW Ensemble (von Mises)		
	FoM (68% CI)	FoM	FoM	FoM	FoM	FoM	FoM	λ	FoM	λ
2.7 keV	0.78 (1.19)	0.76 (1.16)	2.6 (3.3)	0.75 (1.13)	0.74 (1.13)	0.66 (1.01)	0.66 (1.01)	1.76	0.65 (1.0)	1.24
4.5 keV	0.69 (1.05)	0.67 (1.03)	1.5 (1.9)	0.63 (0.94)	0.61 (0.94)	0.56 (0.85)	0.56 (0.85)	1.4	0.55 (0.84)	1.12
6.4 keV	0.58 (0.88)	0.56 (0.86)	1.6 (1.9)	0.50 (0.75)	0.49 (0.74)	0.45 (0.69)	0.45 (0.69)	1.1	0.44 (0.68)	1.02
8.0 keV	0.53 (0.8)	0.51 (0.79)	0.8 (1.1)	0.48 (0.71)	0.46 (0.71)	0.43 (0.66)	0.43 (0.66)	1.08	0.42 (0.65)	1.07
PL2	1.12 (1.72)	1.07 (1.64)	–	1.08 (1.64)	1.07 (1.63)	0.89 (1.36)	0.88 (1.34)	1.85	0.85 (1.29)	1.28
PL1	1.02 (1.56)	0.97 (1.48)	–	0.97 (1.46)	0.95 (1.45)	0.79 (1.2)	0.79 (1.2)	1.69	0.78 (1.18)	1.25

Table 1: Results on energy selected track image datasets, comparing our method with the current state of the art and including an ablation study. Lower FoM is better. PL1 and PL2 are power law datasets with range spanning 2.0 – 8.0keV (PL1 $dN/dE \propto E^{-1}$, and PL2 $dN/dE \propto E^{-2}$). All test datasets have 360 thousand tracks each to enable comparison with Kitaguchi et al. (2019). All methods have $RMSE_\phi \leq 0.5^\circ$. Confidence intervals (CI 68%) are calculated using the non-parametric bootstrap – note these are not the standard errors on FoM values, standard errors are $CI/\sqrt{200}$, except in the case of Kitaguchi et al. (2019). For FoMs only the upper CI bound is necessary, since this represents the worst case signal to noise ratio. IW stands for importance weighted. All of our method results use the Gaussian loss, (5), except for the final column which uses the von Mises loss. All ensembles have $M = 10$ NN members.

258 where $\mathbf{v}_n = (\cos\hat{y}_n, \sin\hat{y}_n)$ and $\mathbf{x} = (\Pi\cos\phi, \Pi\sin\phi)$. By recasting (2) as a convex optimization
 259 problem, we have a guaranteed globally optimal solution for (Π, ϕ) . We can solve (7) quickly and
 260 efficiently using second order Newton methods. In practice we use the robust open source software
 261 IpOpt, §2.4.

262 We also consider a more domain specific, non-Gaussian likelihood function for our loss, (5). We
 263 use the log-likelihood of the von Mises distribution for the NN loss:

$$\text{Loss}(\mathbf{v}, \hat{\mathbf{v}}) = \log(I_0(\hat{\sigma}^{-2})) - \frac{1}{\hat{\sigma}^2} \mathbf{v}^T \hat{\mathbf{v}}, \quad (8)$$

264 where I_0 is the modified Bessel function of the first kind. This is a close approximation of the
 265 wrapped Gaussian on the circle. It is more appropriate than the Gaussian (5) for angular estimates
 266 since it can capture the π periodicity in \hat{y} . For very small $\hat{\sigma}$ this is equivalent to the Gaussian. We
 267 compare the results from both losses in §3.4 and table 1.

268 3.2.1 FIGURE OF MERIT

269 In polarization estimation, we want high recovered $\hat{\Pi}_{100\%}$ (and accurate ϕ) for a known 100%
 270 polarized source ($\Pi = 1$), and low recovered $\hat{\Pi}_{0\%}$ for an unpolarized source ($\Pi = 0$). Since there
 271 is irreducible noise in the tracks, it is impossible for any method to achieve $\hat{\Pi}_{100\%} \sim 1$, so $\hat{\Pi}_{\text{meas}}$
 272 estimates are calibrated to get the final $\hat{\Pi}$ for an unknown source²: $\hat{\Pi} = \hat{\Pi}_{\text{meas}}/\hat{\Pi}_{100\%}$. We define a
 273 figure of merit for polarization estimation:

$$\text{FoM} = 100 \times \hat{\Pi}_{0\%}/\hat{\Pi}_{100\%}. \quad (9)$$

274 We use the FoM to evaluate model performance: a lower FoM means better polarization estimation.
 275 This is effectively a measure of the signal to noise ratio, a simplified extension of the minimum
 276 detectable polarization (MDP) typically defined for X-ray polarization (Weisskopf et al., 2010) that
 277 does not preclude biased estimators. It is evaluated on unseen polarized and unpolarized datasets.
 278 In estimating the FoM, we take the number of tracks $N \sim 360,000$ so we can compare directly to
 279 Kitaguchi et al. (2019). We average the FoM over 200 independent track dataset samples of size N .
 280 We use the FoM as the criterion to select the hyperparameter λ in (2). In this way we can tradeoff
 281 accuracy and bias in our Π, ϕ estimates.

282 3.3 NN TRAINING AND SELECTION

283 Our training dataset consists of 3 million simulated tracks, examples of which are shown in fig. 1.
 284 The track energies uniformly span 1.0 – 9.0keV, IXPE’s most sensitive range and are unpolarized
 285 (uniform track angle distribution). Since we don’t know a priori what energy each track is, we want

² $\hat{\Pi}_{100\%}$ is measured before on a source with the same track energy distribution.

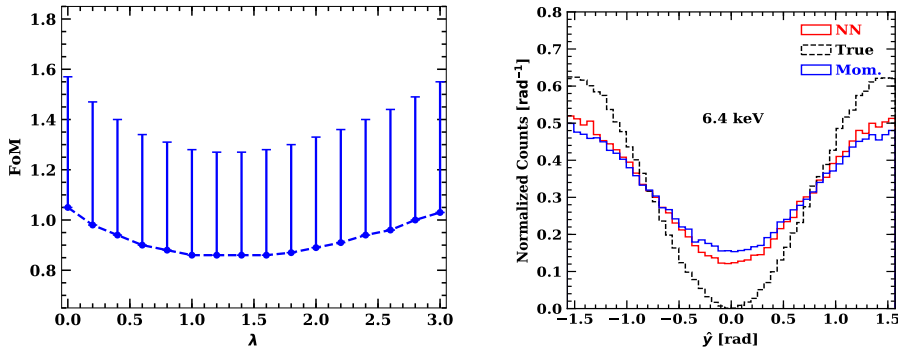


Figure 2: *Left*: FoM as a function of hyperparameter λ for the von Mises ensemble on the PL2 dataset. This method is used to select all of the λ . *Right*: Histogram of \hat{y} predictions for the 6.4KeV polarized dataset, $\Pi = 1$, $\phi = \pi/2$ (ground truth). Black shows the ground truth density, (4), to be estimated. Red and blue show the single NN and standard moments estimates respectively. A single NN can better predict \hat{y} and thus extract more polarization signal, $\Pi_{100\%}$ resulting in a better FoM.

286 NNs that can make predictions for tracks of all energies. This also makes for a more generalizable
 287 system, since some high energy tracks have similar characteristics to lower energy ones. Each track
 288 is labelled with its 2D angle vector \mathbf{v} .

289 We use a ResNet-19 (He et al., 2015) convolutional NN architecture as our base NN. This particular
 290 architecture is large enough to overfit the training set, and trains in a reasonable amount of time.
 291 Before training we preprocess the training data (square track images). We apply pixelwise centering
 292 and rescaling. We use stochastic gradient descent with momentum and a decaying learning rate
 293 starting at $1e - 2$. We choose batch sizes 512, 1024, 2048 (tracks per batch). We trained for 150
 294 epochs, using early stopping to prevent overfitting. We use L_2 -norm regularization 5×10^{-5} . We
 295 train 30 NNs and compare randomly selecting $M = 10$ NNs to selecting $M = 10$ NNs with the top
 296 MSEs on y for an unseen test dataset spanning all energies to make up our final NN ensemble. The
 297 results for both methods are shown in table 1.

298 3.4 RESULTS

299 Table 1 shows the results of our deep ensemble PFDE method alongside the current state of the art
 300 methods. The single CNN method with optimized cuts, developed in (Kitaguchi et al., 2019), pro-
 301 vides significant improvements in $\Pi_{100\%}$ over the moment analysis, but adds bias to the unpolarized
 302 measurement $\Pi_{0\%}$, increasing its FoM and making it a worse method for all energies. We perform
 303 an ablation study over our method, testing a single NN without using weighting when estimating
 304 (Π, ϕ) (i.e. $w_n = 1 \forall n$, (3)), an ensemble of NNs without weighting, a randomly selected ensemble
 305 with weighting, a top MSE selected ensemble with weighting and a von Mises loss weighted en-
 306 semble. We find a single NN without weighting beats the classical moments and moments with cuts
 307 baselines. This result is visualized in the right panel of fig. 3.3 for the 6.4keV dataset: the single NN
 308 shows improved \hat{y} estimates and thus a density that more closely resembles the ground truth. Using
 309 an ensemble of NNs improves this result slightly, but the real power of our method comes with the
 310 importance weights. Our final importance weighted ensemble method, with λ tuned accordingly for
 311 each energy, significantly outperforms the rest, especially in the power law datasets, where there is
 312 a reduction in FoM of almost a factor of 1.5. This shows the power of a simple weighted scheme
 313 over quality cuts in PFDE, it allows our method to take advantage of higher signal ($\Pi_{100\%}$) at higher
 314 energies in the power law datasets. The λ tuning procedure is shown in the left panel of fig.3.3.

315 Comparing a randomly selected ensemble with a top MSE selected ensemble we find the results
 316 are almost identical. Random selection should yield more accurate approximations of the epistemic
 317 uncertainty and thus better weights, while selecting top performing NN on MSE should improve \hat{y}
 318 accuracy. Since the results are identical, but selecting NNs has the potential to bias density estima-
 319 tion, we recommend randomly selecting NNs. We note that, although not included in the table, a
 320 single NN with importance weighting performs only slightly worse than than the weighted ensem-
 321 ble. Since a single NN only produces aleatoric uncertainties, this suggests, as expected, that for a
 322 correctly specified model aleatoric uncertainties dominate epistemic ones. Finally, the von Mises

323 loss shows a small improvement over the simple Gaussian. This is expected, since characterizing
324 the predictive uncertainties by a periodic distribution is more appropriate for the polarimetry ap-
325 plication, but the improvement is small, suggesting that the Gaussian is a robust starting point for
326 many applications. We plan to release further results and more domain specific information for this
327 particular application [*reference deleted to maintain integrity of review process*].

328 3.5 OTHER APPLICATIONS

329 There are numerous application of PFDE with uncertainty in the physical sciences and engineering.
330 In high energy particle physics massive, short-lived particles can be detected by fitting a Cauchy
331 distribution to the frequencies of measured decay states. Raw sensor data from hadronic particle
332 colliders like the LHC are very noisy with variable uncertainty, meaning our PFDE approach to
333 estimate the Cauchy distribution parameters could be very fruitful. This especially true with the
334 widespread current use of deep learning in particle physics (Guest et al., 2018). Our approach is
335 heuristically justified due to the asymptotic efficiency of the maximum likelihood estimator in a
336 Cauchy location model (Cohen Freue, 2007). In manufacturing, GLMs fit to binomial distributions
337 are commonly used to assess product quality, or the probability of a product being defunct. Today,
338 computer vision is used for much of the inspection (Rossol, 1983), making our hybrid PFDE method
339 a potential step forward. These are just a few application examples – our method may be useful for
340 any GLM based method with high dimensional data.

341 4 DISCUSSION

342 We have proposed a supervised learning framework for parametric feature density estimation. Our
343 method uses deep ensembles to predict high dimensional data features, their aleatoric and epistemic
344 uncertainties. We estimate feature density parameters by incorporating both of these uncertainties
345 into an importance weighted maximum likelihood estimate. We include a tuneable weighting hyper-
346 parameter λ , allowing one to control the bias-variance tradeoff for density estimation. Intuitively,
347 in many real feature density estimation problems, some high dimensional data points may be much
348 more informative than others due to complex noise or differing generative distributions. Our method
349 models this explicitly, weighting datapoint features by their predictive uncertainty when estimating
350 density parameters. This avoids throwing away valuable data with quality cuts, yielding improved
351 density estimates. Our method is scaleable to any feature dataset size and is completely flexible for
352 specific domain applications; most NN architectures can be used. We achieve state-of-the-art results
353 over standard deep learning methods and classical algorithms in X-ray polarimetry - a recent open
354 problem in ML. We expect our method would provide similar improvements to a number of PFDE
355 application fields, including high energy particle physics and manufacturing.

356 We perform an ablation study comparing a single NN, a deep ensemble, and various importance
357 weighted deep ensembles. A single NN approach or standard deep ensemble improves slightly
358 on the classical baselines, but importance weighting by predictive uncertainty provides the main
359 improvements to our method. Selecting NNs for the deep ensemble based on quality of density
360 estimation provides no additional gain in performance compared to random selection – since it is
361 possible performance-based NN selection can degrade epistemic uncertainty estimates, we recom-
362 mend randomly selecting NNs for the ensemble. Comparing the Gaussian and von Mises distribution
363 for feature prediction we find the standard Gaussian likelihood (1) an effective and robust approx-
364 imation, although results can potentially be improved for specific applications by choosing a more
365 appropriate distribution over the predictive uncertainties.

366 While our method works well for densities with convex log-likelihoods, non-convex ones will not
367 necessarily yield globally optimal solutions and may be very time consuming to evaluate. *Future*
368 *Work*: Future additions to the method include more complex aleatoric uncertainty modelling. We
369 assume a Gaussian distribution for our feature prediction (1), but for domain applications where
370 there is an expected feature uncertainty, one could use an alternative distribution, or even a mixture
371 density network (Bishop, 1994) for more flexibility. In that case the functional form of weighting
372 would have to be reconsidered. Additionally, finding the optimal weighting function for specific
373 problem applications is likely to yield significant improvements.

374 REFERENCES

- 375 Ronaldo Bellazzini, F. Angelini, Luca Baldini, Alessandro Brez, Enrico Costa, Giuseppe Di
376 Persio, Luca Latronico, M. M. Massai, Nicola Omodei, Luigi Pacciani, Paolo Soffitta,
377 and Gloria Spandre. Novel gaseous x-ray polarimeter: data analysis and simulation.
378 In *Polarimetry in Astronomy*, volume 4843, pp. 383–393. International Society for Op-
379 tics and Photonics, February 2003. doi: 10.1117/12.459381. URL [https://www.
380 spiedigitallibrary.org/conference-proceedings-of-spie/4843/0000/
381 Novel-gaseous-x-ray-polarimeter-data-analysis-and-simulation/
382 10.1117/12.459381.short](https://www.spiedigitallibrary.org/conference-proceedings-of-spie/4843/0000/Novel-gaseous-x-ray-polarimeter-data-analysis-and-simulation/10.1117/12.459381.short).
- 383 Christopher Bishop. Mixture Density Networks. January 1994. URL
384 [https://www.microsoft.com/en-us/research/publication/
385 mixture-density-networks/](https://www.microsoft.com/en-us/research/publication/mixture-density-networks/).
- 386 Ethem F. Can, Aysu Ezen-Can, and Fazli Can. Multilingual Sentiment Analysis: An RNN-Based
387 Framework for Limited Data. *arXiv:1806.04511 [cs]*, June 2018. URL [http://arxiv.org/
388 abs/1806.04511](http://arxiv.org/abs/1806.04511). arXiv: 1806.04511.
- 389 Gabriela V. Cohen Freue. The Pitman estimator of the Cauchy location parameter. *Jour-
390 nal of Statistical Planning and Inference*, 137(6):1900–1913, June 2007. ISSN 0378-3758.
391 doi: 10.1016/j.jspi.2006.05.002. URL [http://www.sciencedirect.com/science/
392 article/pii/S0378375806001285](http://www.sciencedirect.com/science/article/pii/S0378375806001285).
- 393 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep
394 Bidirectional Transformers for Language Understanding. *arXiv:1810.04805 [cs]*, May 2019.
395 URL <http://arxiv.org/abs/1810.04805>. arXiv: 1810.04805.
- 396 B. Efron. Bootstrap Methods: Another Look at the Jackknife. *Annals of Statistics*, 7(1):1–26,
397 January 1979. ISSN 0090-5364, 2168-8966. doi: 10.1214/aos/1176344552. URL [https://
398 projecteuclid.org/euclid.aos/1176344552](https://projecteuclid.org/euclid.aos/1176344552). Publisher: Institute of Mathematical
399 Statistics.
- 400 F. Feroz, M. P. Hobson, and M. Bridges. MultiNest: an efficient and robust Bayesian inference tool
401 for cosmology and particle physics. *Monthly Notices of the Royal Astronomical Society*, 398(4):
402 1601–1614, October 2009. ISSN 00358711, 13652966. doi: 10.1111/j.1365-2966.2009.14548.x.
403 URL <http://arxiv.org/abs/0809.3437>. arXiv: 0809.3437.
- 404 Stanislav Fort, Huiyi Hu, and Balaji Lakshminarayanan. Deep Ensembles: A Loss Landscape Per-
405 spective. *arXiv:1912.02757 [cs, stat]*, December 2019. URL [http://arxiv.org/abs/
406 1912.02757](http://arxiv.org/abs/1912.02757). arXiv: 1912.02757.
- 407 Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist tem-
408 poral classification: labelling unsegmented sequence data with recurrent neural networks. In
409 *Proceedings of the 23rd international conference on Machine learning*, ICML ’06, pp. 369–
410 376, Pittsburgh, Pennsylvania, USA, June 2006. Association for Computing Machinery. ISBN
411 978-1-59593-383-6. doi: 10.1145/1143844.1143891. URL [https://doi.org/10.1145/
412 1143844.1143891](https://doi.org/10.1145/1143844.1143891).
- 413 Dan Guest, Kyle Cranmer, and Daniel Whiteson. Deep Learning and its Application to LHC Physics.
414 *Annual Review of Nuclear and Particle Science*, 68(1):161–181, October 2018. ISSN 0163-8998,
415 1545-4134. doi: 10.1146/annurev-nucl-101917-021019. URL [http://arxiv.org/abs/
416 1806.11484](http://arxiv.org/abs/1806.11484). arXiv: 1806.11484.
- 417 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image
418 Recognition. *arXiv:1512.03385 [cs]*, December 2015. URL [http://arxiv.org/abs/
419 1512.03385](http://arxiv.org/abs/1512.03385). arXiv: 1512.03385.
- 420 Feifang Hu and James V. Zidek. The Weighted Likelihood. *The Canadian Journal of Statistics / La
421 Revue Canadienne de Statistique*, 30(3):347–371, 2002. ISSN 0319-5724. doi: 10.2307/3316141.
422 URL <https://www.jstor.org/stable/3316141>.

- 423 Nikos Karampatziakis and John Langford. Online importance weight aware updates. In *Proceedings*
424 *of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, UAI'11, pp. 392–399,
425 Barcelona, Spain, July 2011. AUAI Press. ISBN 978-0-9749039-7-2.
- 426 Alex Kendall and Yarin Gal. What Uncertainties Do We Need in Bayesian Deep Learning for
427 Computer Vision? In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vish-
428 wanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*,
429 pp. 5574–5584. Curran Associates, Inc., 2017. URL [http://papers.nips.cc/paper/
430 7141-what-uncertainties-do-we-need-in-bayesian-deep-learning-for-computer-vision.
431 pdf](http://papers.nips.cc/paper/7141-what-uncertainties-do-we-need-in-bayesian-deep-learning-for-computer-vision.pdf).
- 432 Takao Kitaguchi, Kevin Black, Teruaki Enoto, Asami Hayato, Joanne E. Hill, Wataru B. Iwakiri,
433 Philip Kaaret, Tsunefumi Mizuno, and Toru Tamagawa. A convolutional neural network ap-
434 proach for reconstructing polarization information of photoelectric X-ray polarimeters. *Nuclear*
435 *Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors*
436 *and Associated Equipment*, 942:162389, October 2019. ISSN 01689002. doi: 10.1016/j.nima.
437 2019.162389. URL <http://arxiv.org/abs/1907.06442>. arXiv: 1907.06442.
- 438 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet Classification with
439 Deep Convolutional Neural Networks. In F. Pereira, C. J. C. Burges, L. Bottou, and
440 K. Q. Weinberger (eds.), *Advances in Neural Information Processing Systems 25*, pp.
441 1097–1105. Curran Associates, Inc., 2012. URL [http://papers.nips.cc/paper/
442 4824-imagenet-classification-with-deep-convolutional-neural-networks.
443 pdf](http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf).
- 444 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive
445 uncertainty estimation using deep ensembles. In *Proceedings of the 31st International Conference*
446 *on Neural Information Processing Systems*, NIPS'17, pp. 6405–6416, Long Beach, California,
447 USA, December 2017. Curran Associates Inc. ISBN 978-1-5108-6096-4.
- 448 Nikita Moriakov, Ashwin Samudre, Michela Negro, Fabian Gieseke, Sydney Otten, and Luc Hen-
449 driks. Inferring astrophysical X-ray polarization with deep learning. *arXiv:2005.08126 [astro-*
450 *ph]*, May 2020. URL <http://arxiv.org/abs/2005.08126>. arXiv: 2005.08126.
- 451 J. A. Nelder and R. W. M. Wedderburn. Generalized Linear Models. *Journal of the Royal Statistical*
452 *Society. Series A (General)*, 135(3):370–384, 1972. ISSN 0035-9238. doi: 10.2307/2344614.
453 URL <https://www.jstor.org/stable/2344614>. Publisher: [Royal Statistical Soci-
454 ety, Wiley].
- 455 Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D. Sculley, Sebastian Nowozin, Joshua V.
456 Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can You Trust Your Model’s Uncertainty?
457 Evaluating Predictive Uncertainty Under Dataset Shift. *arXiv:1906.02530 [cs, stat]*, December
458 2019. URL <http://arxiv.org/abs/1906.02530>. arXiv: 1906.02530.
- 459 George Papamakarios. Neural Density Estimation and Likelihood-free Inference. *arXiv:1910.13233*
460 *[cs, stat]*, October 2019. URL <http://arxiv.org/abs/1910.13233>. arXiv:
461 1910.13233.
- 462 George Papamakarios, Theo Pavlakou, and Iain Murray. Masked Autoregressive Flow for Density
463 Estimation. *arXiv:1705.07057 [cs, stat]*, June 2018. URL [http://arxiv.org/abs/1705.
464 07057](http://arxiv.org/abs/1705.07057). arXiv: 1705.07057.
- 465 Tim Pearce, Mohamed Zaki, and Andy Neely. Bayesian Neural Network Ensembles.
466 *arXiv:1811.12188 [cs, stat]*, November 2018. URL [http://arxiv.org/abs/1811.
467 12188](http://arxiv.org/abs/1811.12188). arXiv: 1811.12188.
- 468 Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global Vectors for Word
469 Representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Lan-*
470 *guage Processing (EMNLP)*, pp. 1532–1543, Doha, Qatar, October 2014. Association for Com-
471 putational Linguistics. doi: 10.3115/v1/D14-1162. URL [https://www.aclweb.org/
472 anthology/D14-1162](https://www.aclweb.org/anthology/D14-1162).

- 473 Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You Only Look Once: Unified,
474 Real-Time Object Detection. *arXiv:1506.02640 [cs]*, June 2015. URL [http://arxiv.org/
475 abs/1506.02640](http://arxiv.org/abs/1506.02640). arXiv: 1506.02640.
- 476 Lothar Rossol. Computer Vision in Industry. In Alan Pugh (ed.), *Robot Vision, International
477 Trends in Manufacturing Technology*, pp. 11–18. Springer, Berlin, Heidelberg, 1983. ISBN 978-
478 3-662-09771-7. doi: 10.1007/978-3-662-09771-7_2. URL [https://doi.org/10.1007/
978-3-662-09771-7_2](https://doi.org/10.1007/
479 978-3-662-09771-7_2).
- 480 Carmelo Sgro. The gas pixel detector on board the IXPE mission. In Oswald H.
481 Siegmund (ed.), *UV, X-Ray, and Gamma-Ray Space Instrumentation for Astron-
482 omy XX*, pp. 16, San Diego, United States, August 2017. SPIE. ISBN 978-1-
483 5106-1251-8 978-1-5106-1252-5. doi: 10.1117/12.2273922. URL [https://www.
484 spiedigitallibrary.org/conference-proceedings-of-spie/10397/
485 2273922/The-gas-pixel-detector-on-board-the-IXPE-mission/10.
486 1117/12.2273922.full](https://www.spiedigitallibrary.org/conference-proceedings-of-spie/10397/2273922/The-gas-pixel-detector-on-board-the-IXPE-mission/10.1117/12.2273922.full).
- 487 C. Sgrò and IXPE Team. The Imaging X-ray Polarimetry Explorer (IXPE). *Nuclear Instru-
488 ments and Methods in Physics Research A*, 936:212–215, August 2019. ISSN 0168-9002. doi:
489 10.1016/j.nima.2018.10.111. URL [http://adsabs.harvard.edu/abs/2019NIMPA.
490 936..212S](http://adsabs.harvard.edu/abs/2019NIMPA.936..212S).
- 491 Constantin Steppa and Tim Lukas Holch. HexagDLY - Processing hexagonally sampled data with
492 CNNs in PyTorch. *SoftwareX*, 9:193–198, January 2019. ISSN 23527110. doi: 10.1016/j.softx.
493 2019.02.010. URL <http://arxiv.org/abs/1903.01814>. arXiv: 1903.01814.
- 494 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
495 Łukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In I. Guyon, U. V. Luxburg,
496 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neu-
497 ral Information Processing Systems 30*, pp. 5998–6008. Curran Associates, Inc., 2017. URL
498 <http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf>.
- 499 Martin Weisskopf. An Overview of X-Ray Polarimetry of Astronomical Sources. *Galaxies*, 6:33,
500 March 2018. doi: 10.3390/galaxies6010033. URL [http://adsabs.harvard.edu/abs/
501 2018Galax...6...33W](http://adsabs.harvard.edu/abs/2018Galax...6...33W).
- 502 Martin C. Weisskopf, Ronald F. Elsner, and Stephen L. O’Dell. On understanding the figures of merit
503 for detection and measurement of x-ray polarization. *arXiv:1006.3711 [astro-ph]*, pp. 77320E,
504 July 2010. doi: 10.1117/12.857357. URL <http://arxiv.org/abs/1006.3711>. arXiv:
505 1006.3711.
- 506 Andreas Wächter and Lorenz T. Biegler. On the implementation of an interior-point filter line-
507 search algorithm for large-scale nonlinear programming. *Mathematical Programming*, 106(1):
508 25–57, March 2006. ISSN 1436-4646. doi: 10.1007/s10107-004-0559-y. URL [https://
509 doi.org/10.1007/s10107-004-0559-y](https://doi.org/10.1007/s10107-004-0559-y).