

The DNNLikelihood: enhancing likelihood distribution with Deep Learning

Andrea Coccaro^a, Maurizio Pierini^b, Luca Silvestrini^{b,c}, and Riccardo Torre^{a,b}

^a *INFN, Sezione di Genova, Via Dodecaneso 33, I-16146 Genova, Italy*

^b *CERN, 1211 Geneva 23, Switzerland*

^c *INFN, Sezione di Roma, P.le A. Moro, 2, I-00185 Roma, Italy*

Abstract

We introduce the DNNLikelihood, a novel framework to easily encode, through Deep Neural Networks (DNN), the full experimental information contained in complicated likelihood functions (LFs). We show how to efficiently parametrise the LF, treated as a multivariate function of parameters and nuisance parameters with high dimensionality, as an interpolating function in the form of a DNN predictor. We do not use any Gaussian approximation or dimensionality reduction, such as marginalisation or profiling over nuisance parameters, so that the full experimental information is retained. The procedure applies to both binned and unbinned LFs, and allows for an efficient distribution to multiple software platforms, e.g. through the framework-independent ONNX model format. The distributed DNNLikelihood can be used for different use cases, such as re-sampling through Markov Chain Monte Carlo techniques, possibly with custom priors, combination with other LFs, when the correlations among parameters are known, and re-interpretation within different statistical approaches, i.e. Bayesian vs frequentist. We discuss the accuracy of our proposal and its relations with other approximation techniques and likelihood distribution frameworks. As an example, we apply our procedure to a pseudo-experiment corresponding to a realistic LHC search for new physics already considered in the literature.

1 Introduction

The Likelihood Function (LF) is the fundamental ingredient of any statistical inference. It encodes the full information on experimental measurements and allows for their interpretation both from a frequentist (e.g. Maximum Likelihood Estimation (MLE)) and a Bayesian (e.g. Maximum a Posteriori (MAP)) perspectives.¹

On top of providing a description of the combined conditional probability distribution of data given a model (or vice versa, when the prior is known, of a model given data), and therefore of the relevant statistical uncertainties, the LF may also encode, through the so-called *nuisance parameters*, the full knowledge of systematic uncertainties and additional constraints (for instance coming from the measurement of fundamental input parameters by other experiments) affecting a given measurement or observation, as for instance discussed in ref. [3].

Current experimental and phenomenological results in fundamental physics and astrophysics typically involve complicated fits with several parameters of interest and hundreds of nuisance parameters. Unfortunately, it is generically considered a hard task to provide all the information encoded in the LF in a practical and reusable way. Therefore, experimental analyses usually deliver only a small fraction of the full information contained in the LF, typically in the form of confidence intervals obtained by profiling the LF on the nuisance parameters (frequentist approach), or in terms of probability intervals obtained by marginalising over nuisance parameters (Bayesian approach), depending on the statistical method used in the analysis. This way of presenting results is very practical, since it can be encoded graphically into simple plots and/or simple tables of expectation values and correlation matrices among observables, effectively making use of the Gaussian approximation. However, such ‘partial’ information can hardly be used to reinterpret a result within a different physics scenario, to combine it with other results, or to project its sensitivity to the future. These tasks, especially outside experimental collaborations, are usually done in a naïve fashion, trying to reconstruct an approximate likelihood for the quantities of interest, employing a Gaussian approximation, assuming full correlation/uncorrelation among parameters, and with little or no control on the effect of systematic uncertainties. Such control on systematic uncertainties could be particularly useful to project the sensitivity of current analyses to future experiments, an exercise particularly relevant in the context of future collider studies [4–6]. One could for instance ask how a certain experimental result would change if a given systematic uncertainty (theoretical or experimental) is reduced by some amount. This kind of question is usually not addressable using only public results for the aforementioned reasons. This and other limitations could of course be overcome if the full LF were available as a function of

¹To be precise, the LF does not contain the full information in the frequentist approach, since the latter does not satisfy the Likelihood Principle (for a detailed comparison of frequentist and Bayesian inference see, for instance, refs. [1,2]). In particular, frequentists should specify assumptions about the experimental setup and about experiments that are not actually performed which are not relevant in the Bayesian approach. Since these assumptions are however usually well spelled in fundamental physics and astrophysics (at least when classical inference is carefully applied), we ignore this issue and assume that the LF encodes the full experimental information.

the observables and of the elementary nuisance parameters, allowing for:

1. the combination of the LF with other LFs involving (a subset of) the same observables and/or nuisance parameters;
2. the reinterpretation of the analysis under different theoretical assumptions (up to issues with unfolding);
3. the reuse of the LF in a different statistical framework;
4. the study of the dependence of the result on the prior knowledge of the observables and/or nuisance parameters.

A big effort has been put in recent years into improving the distribution of information on the experimental LFs, usually in the form of binned histograms in mutually exclusive categories, or, even better, giving information on the covariance matrix between them. An example is given by the Higgs Simplified Template Cross Sections [7]. Giving only information on central values and uncertainties, this approach makes intrinsic use of the Gaussian approximation to the LF without preserving the original information on the nuisance parameters of each given analysis. A further step has been taken in refs. [8–10] where a simplified parameterisation in terms of a set of “effective” nuisance parameters was proposed, with the aim of catching the main features of the true distribution of the observables, up to the third moment.² This is a very practical and effective solution, sufficiently accurate for many use cases. On the other hand, its underlying approximations come short whenever the dependence of the LF on the original nuisance parameters is needed, and with highly non-Gaussian (e.g. multi-modal) LFs.

Recently, the ATLAS collaboration has taken a major step forward, releasing the full experimental LF of an analysis [12] on HEPData [13] through the HistFactory framework [14], with the format presented in ref. [15]. Before this, the release of the full experimental likelihood was advocated several times as a fundamental step forward for the HEP community (see e.g. the panel discussion in ref. [16]), but so far it was not followed up with a concrete commitment. The lack of a concrete effort in this direction was usually attributed to technical difficulties and indeed, this fact was our main motivation when we initiated the work presented in this paper.

In this work, we propose to present the full LF, as used by experimental collaborations to produce the results of their analyses, in the form of a suitably trained Deep Neural Network (DNN), which is able to reproduce the original LF as a function of physical and nuisance parameters with the accuracy required to allow for the four aforementioned tasks.

The DNNlikelihood approach offers at least two remarkable practical advantages. First, it does not make underlying assumptions on the structure of the LF (e.g. binned vs unbinned), extending to use cases that might be problematic for currently available alternative solutions. For instance, there are extremely relevant analyses that are carried out using an unbinned LF,

²This approach is similar to the one already proposed in ref. [11].

notably some Higgs study in the four-leptons golden decay mode [17] and the majority of the analyses carried out at B-physics experiments. Second, the use of a DNN does not impose on the user any specific software choice. Neural Networks are extremely portable across multiple software environments (e.g. C++, PYTHON, Matlab, R, or Mathematica) through the ONNX format [18]. This aspect could be important whenever different experiments make different choices in terms of how to distribute the likelihood.

In this respect, we believe that the use of the DNNLikelihood could be relevant even in a future scenario in which every major experiment has followed the remarkable example set by ATLAS. For instance, it could be useful to overcome some technical difficulty with specific classes of analyses. Or, it could be seen as a further step to take in order to import a distributed likelihood into a different software environment. In addition, the DNNLikelihood could be used in other contexts, e.g. to distribute the outcome of phenomenological analyses involving multi-dimensional fits such as the Unitarity Triangle Analysis [19–23], the fit of electroweak precision data and Higgs signal strengths [24–28], etc.

There are two main challenges associated to our proposed strategy: on one hand, in order to design a supervised learning technique, an accurate sampling of the LF is needed for the training of the DNNLikelihood. On the other hand a (complicated) interpolation problem should be solved with an accuracy that ensures a real preservation of all the required information on the original probability distribution.

The first problem, i.e. the LF sampling, is generally easy to solve when the LF is a relatively simple function which can be fastly evaluated in each point in the parameter space. In this case, Markov Chain Monte Carlo (MCMC) techniques [29] are usually sufficient to get dense enough samplings of the function, even for very high dimensionality, in a reasonable time. However the problem may quickly become intractable with these techniques when the LF is more complicated, and takes much longer time to be evaluated. This is typically the case when the sampling workflow requires the simulation of a data sample, including computationally costly corrections (e.g. radiative corrections) and/or a simulation of the detector response, e.g. through Geant4 [30]. In these cases, evaluating the LF for a point of the parameter space may require $\mathcal{O}(\text{minutes})$ to go through the full chain of generation, simulation, reconstruction, event selection, and likelihood evaluation, making the LF sampling with standard MCMC techniques impractical. To overcome this difficulty, several ideas have recently been proposed, inspired by Bayesian optimisation and Gaussian processes, known as Active Learning (see, for instance, refs. [31, 32] and references therein). These techniques, though less robust than MCMC ones, allow for a very “query efficient” sampling, i.e. a sampling that requires the smallest possible number of evaluations of the full LF. Active Learning applies machine learning techniques to design the proposal function of the sampling points and can be shown to be much more query efficient than standard MCMC techniques. Another possibility would be to employ deep learning and MCMC techniques together in a way similar to Active Learning, but inheriting some of the nice properties of MCMC. We defer a discussion of this new idea to a forthcoming publication [33], while in this work we focus on the second of the aforementioned tasks: we assume that an accurate sampling of the LF is available and design a technique to encode it and distribute it through

DNNs.

A JUPYTER notebook and the PYTHON source files which allow to reproduce all results presented in this paper are available on GitHub [🔗](#). A dedicated PYTHON package allowing to sample LFs and to build, optimize, train, and store the corresponding DNNLikelihoods is in preparation. This will allow not only allow to construct and distribute DNNLikelihoods, but also to use them for inference both in the Bayesian and frequentist frameworks.

This paper is organized as follows. In Section 2 we discuss the issue of interpolating the LF from a Bayesian and frequentist perspective and set up the procedure for producing suitable training datasets. In Section 3 we describe a benchmark example, consisting in the realistic LHC-like New Physics (NP) search proposed in ref. [9]. In Section 4 we show our proposal at work for this benchmark example, whose LF depends on one physical parameter, the signal strength μ , and 94 nuisance parameters. Finally, in Section 5 we conclude and discuss some interesting ideas for future studies.

2 Interpolation of the Likelihood Function

The problem of fitting high-dimensional multivariate functions is a classical interpolation problem, and it is nowadays widely known that DNNs provide the best solution to it. Nevertheless, the choice of the loss function to minimise and of the metrics to quantify the performance, i.e. of the “distance” between the fitted and the true function, depends crucially on the nature of the function and its properties. The LF is a special function, since it represents a probability distribution. As such, it corresponds to the integration measure over the probability of a given set of random variables. The interesting regions of the LF are twofold. From a frequentist perspective the knowledge of the local maxima, and of the global maximum, is needed. This requires a good knowledge of the LF in regions of the parameter space with high probability (large likelihood), and, especially for high dimensionality, very low probability mass (very small prior volume). These regions are therefore very hard to populate via sampling techniques [34] and give tiny contributions to the LF integral, the latter being increasingly dominated by the “tails” of the multidimensional distribution as the number of dimensions grows. From a Bayesian perspective the expectation values of observables or parameters, which can be computed through integrals over the probability measure, are instead of interest. In this case one needs to accurately know regions of very small probabilities, which however correspond to large prior volumes, and could give large contributions to the integrals.

Let us argue what is a good “distance” to minimise to achieve both of the aforementioned goals, i.e. to know the function equally well in the tails and close to the local maxima. Starting from the view of the LF as a probability measure (Bayesian perspective), the quantity that one is interested in minimising is the difference between the expectation values of observables computed using the true probability distribution and the fitted one.

For instance, in a Bayesian analysis one may be interested in using the probability density $\mathcal{P} = \mathcal{L} \times \Pi$, where \mathcal{L} denotes the likelihood and Π the prior, to estimate expectation values

as

$$\mathbf{E}_{\mathcal{P}(\mathbf{x})}[f(\mathbf{x})] = \int f(\mathbf{x})d\mathcal{P}(\mathbf{x}) = \int f(\mathbf{x})\mathcal{P}(\mathbf{x})d\mathbf{x}, \quad (1)$$

where the probability measure is $d\mathcal{P}(\mathbf{x}) = \mathcal{P}(\mathbf{x})d\mathbf{x}$, and we collectively denoted by the n -dimensional vector \mathbf{x} the parameters on which f and \mathcal{P} depend, treating on the same footing the nuisance parameters and the parameters of interest. Let us assume now that the solution to our interpolation problem provides a predicted pdf $\mathcal{P}_P(\mathbf{x})$, leading to an estimated expectation value

$$\mathbf{E}_{\mathcal{P}_P(\mathbf{x})}[f(\mathbf{x})] = \int f(\mathbf{x})\mathcal{P}_P(\mathbf{x})d\mathbf{x}. \quad (2)$$

This can be rewritten, by defining the ratio $r(\mathbf{x}) \equiv \mathcal{P}_P(\mathbf{x})/\mathcal{P}(\mathbf{x}) = \mathcal{L}_P(\mathbf{x})/\mathcal{L}(\mathbf{x})$, as

$$\mathbf{E}_{\mathcal{P}_P(\mathbf{x})}[f(\mathbf{x})] = \int f(\mathbf{x})r(\mathbf{x})\mathcal{P}(\mathbf{x})d\mathbf{x}, \quad (3)$$

so that the absolute error in the evaluation of the expectation value is given by

$$|\mathbf{E}_{\mathcal{P}(\mathbf{x})}[f(\mathbf{x})] - \mathbf{E}_{\mathcal{P}_P(\mathbf{x})}[f(\mathbf{x})]| = \left| \int f(\mathbf{x})(1 - r(\mathbf{x}))\mathcal{P}(\mathbf{x})d\mathbf{x} \right|. \quad (4)$$

For a finite sample of points \mathbf{x}_i , with $i = 1, \dots, N$, the integrals are replaced by sums and eq. (4) becomes

$$|\mathbf{E}_{\mathcal{P}(\mathbf{x})}[f(\mathbf{x})] - \mathbf{E}_{\mathcal{P}_P(\mathbf{x})}[f(\mathbf{x})]| = \left| \frac{1}{N} \sum_{\mathbf{x}_i|_{\mathcal{U}(\mathbf{x})}} f(\mathbf{x}_i)(1 - r(\mathbf{x}_i))\mathcal{F}(\mathbf{x}_i) \right|. \quad (5)$$

Here, the probability density function $\mathcal{P}(\mathbf{x})$ has been replaced with $\mathcal{F}(\mathbf{x}_i)$, the frequencies with which each of the \mathbf{x}_i occurs, normalised such that $\sum_{\mathbf{x}_i|_{\mathcal{U}(\mathbf{x})}} \mathcal{F}(\mathbf{x}_i) = N$, the notation $\mathbf{x}_i|_{\mathcal{U}(\mathbf{x})}$ indicates that the \mathbf{x}_i are drawn from a uniform distribution, and the $1/N$ factor ensures the proper normalisation of probabilities. This sum is very inefficient to calculate when the probability distribution $\mathcal{P}(\mathbf{x})$ varies rapidly in the parameter space, i.e. deviates strongly from a uniform distribution, since most of the \mathbf{x}_i points drawn from the uniform distribution will correspond to very small probabilities, giving negligible contributions to the sum. An example is given by multivariate normal distributions, where, increasing the dimensionality, tails become more and more relevant (see Appendix A). A more efficient way of computing the sum is given by directly sampling the \mathbf{x}_i points from the probability distribution $\mathcal{P}(\mathbf{x})$, so that eq. (5) can be rewritten as

$$|\mathbf{E}_{\mathcal{P}(\mathbf{x})}[f(\mathbf{x})] - \mathbf{E}_{\mathcal{P}_P(\mathbf{x})}[f(\mathbf{x})]| = \left| \frac{1}{N} \sum_{\mathbf{x}_i|\mathcal{P}(\mathbf{x})} f(\mathbf{x}_i)(1 - r(\mathbf{x}_i)) \right|. \quad (6)$$

This expression clarifies the aforementioned importance of being able to sample points from the probability distribution \mathcal{P} to efficiently discretize the integrals and compute expectation values. The minimum of this function for any $f(\mathbf{x})$ is in $r(\mathbf{x}_i) = 1$, which, in turn, implies $\mathcal{L}(\mathbf{x}_i) = \mathcal{L}_{\mathcal{P}}(\mathbf{x}_i)$. This suggests that an estimate of the performance of the interpolated likelihood is given by the Mean Percentage Error (MPE)

$$\text{MPE}_{\mathcal{L}} = \frac{1}{N} \sum_{\mathbf{x}_i | \mathcal{P}(\mathbf{x})} \left(1 - \frac{\mathcal{L}_{\mathcal{P}}(\mathbf{x}_i)}{\mathcal{L}(\mathbf{x}_i)} \right) = \frac{1}{N} \sum_{\mathbf{x}_i | \mathcal{P}(\mathbf{x})} (1 - r(\mathbf{x}_i)) . \quad (7)$$

Technically, formulating the interpolation problem on the LF itself introduces the difficulty of having to fit the function over several orders of magnitude, which leads to numerical instabilities. For this reason it is much more convenient to formulate the problem using the natural logarithm of the LF, the so-called log-likelihood $\log \mathcal{L}$. Let us see how the error on the log-likelihood propagates to the actual likelihood. Consider the Mean Error (ME) on the log-likelihood

$$\text{ME}_{\log \mathcal{L}} = \frac{1}{N} \sum_{\mathbf{x}_i | \mathcal{P}(\mathbf{x})} (\log \mathcal{L}(\mathbf{x}_i) - \log \mathcal{L}_{\mathcal{P}}(\mathbf{x}_i)) = \frac{1}{N} \sum_{\mathbf{x}_i | \mathcal{P}(\mathbf{x})} \log r(\mathbf{x}_i) . \quad (8)$$

The last logarithm can be expanded for $r(\mathbf{x}_i) \sim 1$ to give

$$\text{ME}_{\log \mathcal{L}} \approx \frac{1}{N} \sum_{\mathbf{x}_i | \mathcal{P}(\mathbf{x})} (1 - r(\mathbf{x}_i)) = \text{MPE}_{\mathcal{L}} . \quad (9)$$

It is interesting to notice that $\text{ME}_{\log \mathcal{L}}$ defined in eq. (8) corresponds to the Kullback–Leibler divergence [35], or relative entropy, between \mathcal{P} and $\mathcal{P}_{\mathcal{P}}$:

$$D_{\text{KL}} = \int \log \left(\frac{\mathcal{P}(\mathbf{x})}{\mathcal{P}_{\mathcal{P}}(\mathbf{x})} \right) \mathcal{P}(\mathbf{x}) \, d\mathbf{x} = \frac{1}{N} \sum_{\mathbf{x}_i | \mathcal{P}(\mathbf{x})} \log r(\mathbf{x}_i) = \text{ME}_{\log \mathcal{L}} . \quad (10)$$

While eq. (10) confirms that small values of $D_{\text{KL}} = \text{ME}_{\log \mathcal{L}} \sim \text{MPE}_{\mathcal{L}}$ correspond to a good performance of the interpolation, D_{KL} , as well as ME and MPE, do not satisfy the triangular inequality and therefore cannot be directly optimised for the purpose of training and evaluating a DNN. Eq. (8) suggests however that the Mean Absolute Error (MAE) or the Mean Square Error (MSE) on $\log \mathcal{L}$ should be suitable losses for the DNN training: we explicitly checked that this is indeed the case, with MSE performing slightly better for well-known reasons.

Finally, in the frequentist approach, the LF can be treated just as any function in a regression (or interpolation) problem, and, as we will see, the MSE provides a good choice for the loss function.

2.1 Evaluation metrics

We have argued above that the MAE or MSE on $\log \mathcal{L}(\mathbf{x}_i)$ are the most suitable loss functions to train our DNN for interpolating the LF on the sample \mathbf{x}_i . We are then left with the question of measuring the performance of our interpolation from the statistical point of view. In addition to D_{KL} , several quantities can be computed to quantify the performance of the predictor. First of all, we perform a Kolmogorov-Smirnov (K-S) two-sample test [36,37] on all the marginalised one-dimensional distributions obtained using \mathcal{P} and \mathcal{P}_P . In the hypothesis that both distributions are drawn from the same pdf, the p -value should be distributed uniformly in the interval $[0, 1]$. Therefore, the median of the distribution of p -values of the one-dimensional K-S tests is a representative single number which allows to evaluate the performance of the model. We also compute the error on the width of Highest Posterior Density Intervals (HPDI) PI^i for the marginalised one-dimensional distribution of the i -th parameter, $E_{\text{PI}}^i = |\text{PI}^i - \text{PI}_P^i|$, as well as the relative error on the median of each marginalised distribution. From a frequentist point of view, we are interested in reproducing as precisely as possible the test statistics used in classical inference. In this case we evaluate the model looking at the mean error on the test statistics t_μ , that is the likelihood ratio profiled over nuisance parameters.

To simplify the presentation of the results, we choose the best models according to the median K-S p -value, when considering bayesian inference, and the mean error on the t_μ test-statistics, when considering frequentist inference. These quantities are compared for all the different models on an identical test set statistically independent from both the training and validation sets used for the hyperparameter optimisation.

2.2 Learning from imbalanced data

The loss functions we discussed above are all averages over all samples and, as such, will lead to a better learning in regions that are well represented in the training set and to a less good learning in regions that are under-represented. On the other hand, an unbiased sampling of the LF will populate much more regions corresponding to a large probability mass than regions of large LF. Especially in large dimensionality, it is prohibitive, in terms of the needed number of samples, to make a proper unbiased sampling of the LF, i.e. converging to the underlying probability distribution, while still covering the large LF region with enough statistics. In this respect, learning a multi-dimensional LF raises the issue of learning from highly imbalanced data. This issue is extensively studied in the ML literature for classification problems, but has gathered much less attention from the regression point of view [38–40].

There are two main approaches in the case of regression, both resulting in assigning different weights to different examples. In the first approach, the training set is modified by oversampling and/or undersampling different regions (possibly together with noise) to counteract low/high population of examples, while in the second approach the loss function is modified to weigh more/less regions with less/more examples. In the case where the theoretical underlying distribution of the target variable is (at least approximately) known,

as in our case, either of these two procedures can be applied by assigning weights that are proportional to the inverse frequency of each example in the population. This approach, applied for instance by adding weights to a linear loss function, would really weigh each example equally, which may not be exactly what we need. Moreover, in the case of large dimensionality, the interesting region close to the maximum would be completely absent from the sampling, making any reweighting irrelevant. In this paper we therefore apply an approach belonging to the first class mentioned above, consisting in sampling the LF in the regions of interest and in constructing a training sample that effectively weighs the most interesting regions. As we clarify in Section 3, this procedure consists in building three samples: an unbiased sample, a biased sample and a mixed one. Training data will be extracted from the latter sample. Let us briefly describe the three:

- **Unbiased sample:** a sample that has converged as accurately as possible to the true probability distribution. Notice that this sample is the only one which allows posterior inference in a Bayesian perspective, but would generally fail in making frequentist inference [34].
- **Biased sample:** a sample concentrated around the region of maximum likelihood. It is obtained by biasing the sampler in the region of large LF, only allowing for small moves around the maximum. Tuning this sample, targeted to a frequentist MLE perspective, raises the issue of coverage, that we discuss in Section 3. One has to keep in mind that the region of the LF that needs to be well known, i.e. around the maximum, is related to the coverage of the frequentist analysis being carried out.
- **Mixed sample:** this sample is built by enriching the unbiased sample with the biased one, in the region of large values of the LF. This is a tuning procedure, since, depending on the number of available samples and the statistics needed for training the DNN, this sample needs to be constructed for reproducing at best the results of the given analysis of interest both in a Bayesian and frequentist inference framework.

Some considerations are in order. The unbiased sample is enough if one wants to produce a DNNLikelihood to be used only for Bayesian inference. As we show later, this does not require a complicated tuning of hyperparameters (at least in the example we consider) and reaches very good performances, evaluated with the metrics that we discussed above, already with a relatively small statistics in the training sample (considering the high dimensionality). The situation complicates a bit when one wants to be able to also make frequentist inference using the same DNNLikelihood. In this case the mixed sample (and therefore the biased one) is needed, and more tuning of the network as well as more samples in the training set are required. In order to optimise the predictions made by the DNNLikelihood close to the maximum (or in other words to get an accurate estimate of the t_μ test statistics), allowing for a reliable frequentist inference, we average over a few models with the same architecture, but trained with different randomly generated training sets, therefore employing ensemble learning. This averaging can be done both at the level of the predicted quantities, or at the level of the DNNs. As an example we show the results obtained by averaging the values of

the t_μ test-statistics obtained with a few models trained with different training sets. We have obtained very similar results also by stacking these models together, training a small neural network to optimise their combined prediction, and using this neural network to predict t_μ .

The final issue we have to address when training with the mixed sample, which is biased by construction, is to ensure that the DNNLikelihood can still produce accurate enough Bayesian posterior estimates. This is actually guaranteed by the fact that a regression (or interpolation) problem, contrary to a classification one, is insensitive to the distribution in the target variable, since the output is not conditioned on such probability distribution. This, as can be clearly seen from the results presented in Section 3, is a crucial ingredient for our procedure to be useful, and leads to the main result of our approach: a DNNLikelihood trained with the mixed sample can be used to perform a new MCMC that converges to the underlying distribution, forgetting the biased nature of the training set.

In the next Section we give a thorough example of the procedure discussed here in the case of a prototype LHC-like search for NP corresponding to a 95-dimensional LF.

3 A realistic LHC-like NP search

In this section we introduce the prototype LHC-like NP search presented in ref. [9], which we take as a representative example to illustrate how to train the DNNLikelihood. We refer the reader to ref. [9] for a detailed discussion of this setup and repeat here only the information that is strictly necessary to follow our analysis.

The toy experiment consists in a typical “shape analysis” in a given distribution aimed at extracting information on a possible NP signal from the Standard Model (SM) background. The measurement is divided in three different event categories, containing 30 bins each. The signal is characterized by a single “signal-strength” parameter μ and the uncertainty on the signal is neglected.³ All uncertainties affecting the background are parametrised in terms of nuisance parameters, which may be divided into three categories:

1. fully uncorrelated uncertainties in each bin: they correspond to a nuisance parameter for each bin $\delta_{\text{MC},i}$, with uncorrelated priors, parametrising the uncertainty due to the limited Monte Carlo statistics, or statistics in a control region, used to estimate the number of background events in each bin.
2. fully correlated uncertainties in each bin: they correspond to a single nuisance parameter for each source of uncertainty affecting in a correlated way all bins in the distribution. In this toy experiment, such sources of uncertainty are the modeling of the Initial State Radiation and the Jet Energy Scale, parametrised respectively by the nuisance parameters δ_{ISR} and δ_{JES} .
3. uncertainties on the overall normalisation (correlated among event categories): they correspond to the previous two nuisance parameters δ_{ISR} and δ_{JES} , that, on top of

³This approximation is done in ref. [9] to simplify the discussion, but it is not a necessary ingredient, neither there nor here.

affecting the shape, also affect the overall normalisation in the different categories, plus two typical experimental uncertainties, that only affect the normalisation, given by a veto efficiency and a scale-factor appearing in the simulation, parametrised respectively by δ_{LV} and δ_{RC} .

In summary, the LF depends on one physical parameter μ and 94 nuisance parameters, that we collectively indicate with the vector $\boldsymbol{\delta}$, whose components are defined by $\delta_i = \delta_{MC,i}$ for $i = 1, \dots, 90$, $\delta_{91} = \delta_{ISR}$, $\delta_{92} = \delta_{JES}$, $\delta_{93} = \delta_{LV}$, $\delta_{94} = \delta_{RC}$.

The full model likelihood can be written as⁴

$$\mathcal{L}(\mu, \boldsymbol{\delta}) = \prod_{I=1}^P \Pr(n_I^{\text{obs}} | n_I(\mu, \boldsymbol{\delta})) \pi(\boldsymbol{\delta}), \quad (11)$$

where the product runs over all bins I . The number of expected events in each bin is given by $n_I(\mu, \boldsymbol{\delta}) = n_{s,I}(\mu) + n_{b,I}(\boldsymbol{\delta})$, and the probability distributions are given by Poisson distributions in each bin

$$\Pr(n_I^{\text{obs}} | n_I) = \frac{(n_I)^{n_I^{\text{obs}}} e^{-n_I}}{n_I^{\text{obs}}!}. \quad (12)$$

In this toy LF, the number of background events in each bin $n_{b,I}(\boldsymbol{\delta})$ is known analytically as a function of the nuisance parameters, through various numerical parameters that interpolate the effect of systematic uncertainties. The parametrisation of $n_{b,I}(\boldsymbol{\delta})$ is such that the nuisance parameters $\boldsymbol{\delta}$ are normally distributed with vanishing vector mean and identity covariance matrix

$$\pi(\boldsymbol{\delta}) = \frac{e^{-\frac{1}{2}|\boldsymbol{\delta}|^2}}{(2\pi)^{\frac{\dim(\boldsymbol{\delta})}{2}}}. \quad (13)$$

Moreover, due to the interpolations involved in the parametrisation of the nuisance parameters, in order to ensure positive probabilities, the $\boldsymbol{\delta}$ s are only allowed to take values in the range $[-5, 5]$.

In our approach, we are interested in setting up a supervised learning problem to learn the LF as a function of the parameters. Independently of the statistical perspective, i.e. whether the parameters are treated as random variables or just variables, we need to choose

⁴There is a difference in the interpretation of this formula in the frequentist and Bayesian approaches: in a frequentist approach, the nuisance parameter distributions $\pi(\boldsymbol{\delta})$ do not constitute a prior, but should instead be considered as the likelihood of the nuisance parameters arising from other (auxiliary) measurements [41]. In this perspective, since the product of two likelihoods is still a likelihood, the right hand side of eq. (11) is the full likelihood. On the contrary, in a Bayesian perspective, the full likelihood is given by the product of probabilities in the right hand side of eq. (11), while the distributions $\pi(\boldsymbol{\delta})$ parametrise the *prior* knowledge of the nuisance parameters. Therefore, in this case, according to Bayes' theorem the right hand side of the equation should not be interpreted as the likelihood $\mathcal{P}(\text{data}|\text{pars})$, but as the full posterior probability $\mathcal{P}(\text{pars}|\text{data})$, up to a normalisation given by the Bayesian evidence $\mathcal{P}(\text{data})$. Despite this difference, in order to carry on a unified approach without complicating formulæ too much, we abuse the notation and denote with $\mathcal{L}(\mu, \boldsymbol{\delta})$ the frequentist likelihood and the Bayesian posterior distribution, since these are the two central objects from which frequentist and Bayesian inference are carried out, respectively.

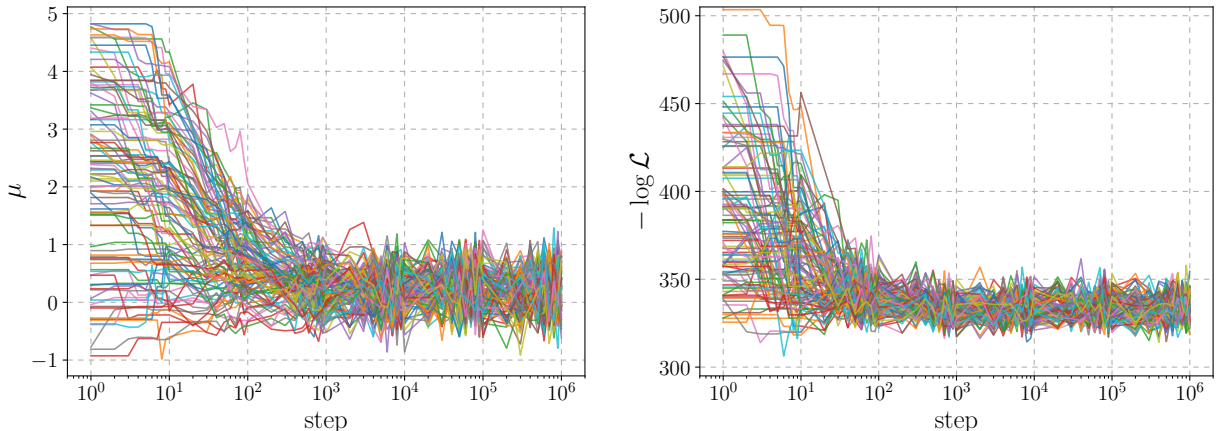


Figure 1: Evolution of the chains in an EMCEE3 sampling of the LF in eq. (11) with 10^3 walkers and 10^6 steps using the STRETCHMOVE algorithm with $a = 1.3$. The plots show the explored values of the parameter μ (left) and of minus log-likelihood $-\log \mathcal{L}$ (right) versus the number of steps for a random subset of 10^2 of the 10^3 chains. For visualization purposes, values in the plots are computed only for numbers of steps included in the set $\{a \cdot 10^b\}$ with $a \in [1, 9]$ and $b \in [0, 6]$.

some values to evaluate the LF. For the nuisance parameters the function $\pi(\boldsymbol{\delta})$ already tells us how to choose these points, since it implicitly treats the nuisance parameters as random variables distributed according to this probability distribution. For the model parameters, in this case only μ , we have to decide how to generate points, independently of the stochastic nature of the parameter itself. In the case of this toy example, since we expect μ to be relatively “small” and most probably positive, we generate μ values according to a uniform probability distribution in the interval $[-1, 5]$. This could be considered as the prior on the stochastic variable μ in a Bayesian perspective, while it is just a scan in the parameter space of μ in the frequentist one.⁵ Notice that we allow for small negative values of μ . Whenever the NP contribution comes from the on-shell production of some new physics, this assumption is not consistent. However, the “signal” may come, in an Effective Field Theory (EFT) perspective, from the interference of the SM background with higher dimensional operators. This interference could be negative depending on the sign of the corresponding Wilson coefficient, and motivates our choice to allow for negative values of μ in our scan.

3.1 Sampling the full likelihood

To obtain the three samples discussed in Section 2.2 from the full model LF in eq. (11) we used the EMCEE3 Python package [42], which implements the Affine Invariant (AI) MCMC

⁵Each different choice of μ corresponds, in the frequentist approach, to a different theoretical hypothesis. This raises the issue of generating pseudo-experiments for each different value of μ , that we discuss further in Section 3.1.2 and Appendix B.

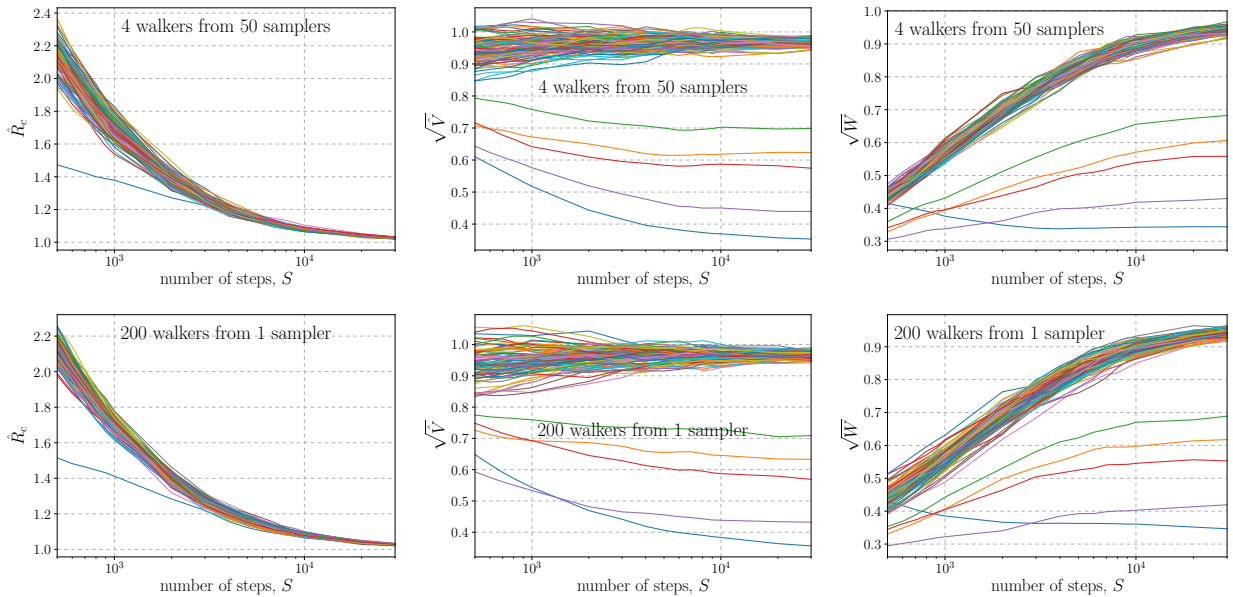


Figure 2: Upper panel: Gelman, Rubin and Brooks \hat{R}_c , $\sqrt{\hat{V}}$, and $\sqrt{\bar{W}}$ parameters computed from an ensemble of 200 walkers made by joining together 4 samples extracted (randomly) from 50 samplers of 200 walkers each. **Lower panel:** Same plots made using 200 walkers from a single sampler.

Ensemble Sampler [43]. We proceeded as follows:

1. Unbiased sample S_1

In the first sampling, the values of the proposals have been updated using the default STRETCHMOVE algorithm implemented in EMCEE3, which updates all values of the parameters (95 in our case) at a time. The default value of the only free parameter of this algorithm $a = 2$ delivered a slightly too low acceptance rate $\epsilon \approx 0.12$. We have therefore set $a = 1.3$, which delivers a better acceptance rate of about 0.36. Walkers⁶ have been initialised randomly according to the prior distribution of the parameters. The algorithm efficiently achieves convergence to the true target distribution, but, given the large dimensionality, hardly explores large values of the LF.

In Figure 1 we show the evolution of the walkers for the parameter μ (left) together with the corresponding values of $-\log \mathcal{L}$ (right) for an illustrative set of 100 walkers. From these figures a reasonable convergence seems to arise already after roughly 10^3 steps, which gives an empirical estimate of the autocorrelation of samples within each walker.

Notice that, in the case of ensemble sampling algorithms, the usual Gelman, Rubin

⁶Walkers are the analog of chains for ensemble sampling methods [43]. In the following, we interchangeably use the words “chains” and “walkers” to refer to the same object.

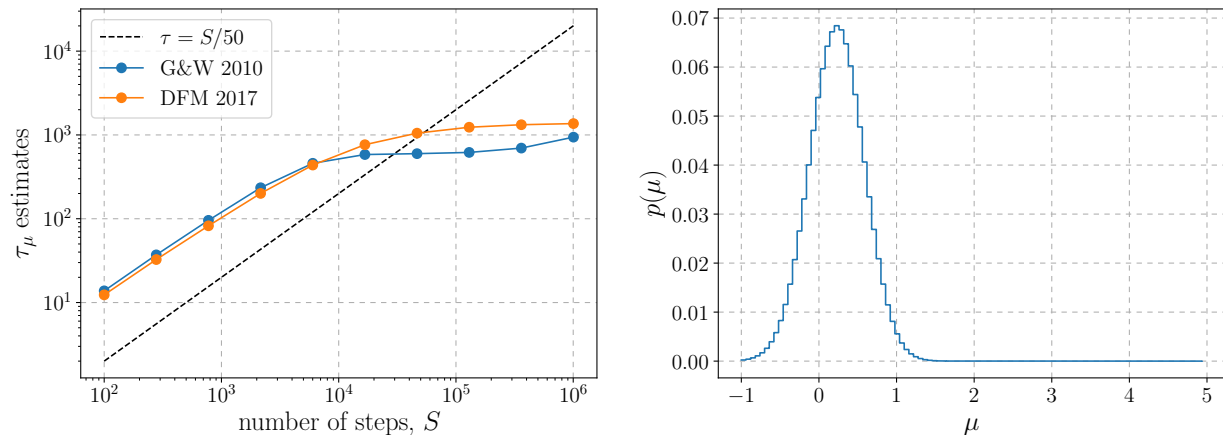


Figure 3: Normalised histogram (right) and estimate of the autocorrelation time τ_μ (estimated using both the original method proposed in ref. [43] and the alternative one discussed by the EMCEE3 authors [47, 48]) as a function of the number of samples (left) for the parameter μ .

and Brooks statistics, usually denoted as \hat{R}_c [44, 45], is not expected to be a robust tool to assess convergence, due to the correlation between different walkers. However, in order to reduce this correlation, one could consider a number of independent samplers, extract a few walkers from each run, and compute \hat{R}_c for this set [46]. Considering the aforementioned empirical estimate of the number of steps for convergence, i.e. roughly few 10^3 , we have run 50 independent samplers for a larger number of steps ($3 \cdot 10^4$), extracted randomly 4 chains from each, joined them together, and computed \hat{R}_c for this set. This is shown in the upper-left plot of Figure 2. With a requirement of $\hat{R}_c < 1.2$ [45] we see that chains have already converged after around $5 \cdot 10^3$ steps, which is roughly what we empirically estimated looking at the chains evolution in Figure 1. An even more robust requirement for convergence is given by $\hat{R}_c < 1.1$, together with stabilized evolution of both variances \hat{V} and W [45]. In the center and right plots of Figure 2 we show this evolution, from which we see that convergence has robustly occurred after $2 - 3 \cdot 10^4$ steps. In order to check the statement that the \hat{R}_c metric cannot be directly applied to ensemble sampling techniques, we have performed the same analysis using 200 walkers from a single sampler. The result is shown in the lower panels of Figure 2. As it can be seen comparing the upper and lower plots, we would have got to the exact same conclusions about convergence, showing that the correlation of walkers does not, at least in this case, affect diagnostics of convergence based on the \hat{R}_c measure.

An alternative and pretty general way to diagnose MCMC sampling is the autocorrelation of chains, and in particular the Integrated Autocorrelation Time (IAT). This quantity represents the average number of steps between two independent samples in the chain. For unimodal distributions, one can generally assume that after a few IAT

the chain forgot where it started and converged to generating samples distributed according to the underlying target distribution. There are more difficulties in the case of multimodal distributions, which are however shared by most of the MCMC convergence diagnostics. We do not enter here in such a discussion, and refer the reader to the overview presented in ref. [49]. An exact calculation of the IAT for large chains is computationally prohibitive, but there are several algorithms to construct estimators of this quantity. The EMCEE3 package comes with tools that implement some of these algorithms, which we have used to study our sampling [47, 48]. To obtain a reasonable estimate of the IAT τ , one needs enough samples, a reasonable empirical estimate of which, that works well also in our case, is at least 50τ . An illustration of this, for the parameter μ , is given in the left panel of Figure 3, where we show, for a sampler with 10^3 chains and 10^6 steps, the IAT estimated after different numbers of steps with two different algorithms, “G&W 2010” and “DFM 2017” (see refs. [47, 48] for details). It is clear from the plot that the estimate becomes flat, and therefore converges to the correct value of the IAT, roughly when the estimate curves cross the empirical value of 50τ (this is an order of magnitude estimate, and obviously, the larger the number of steps, the better the estimate of τ). The best estimate that we get for this sampling for the parameter μ is obtained with 10^6 steps using the “DFM 2017” method and gives $\bar{\tau} \approx 1366$, confirming the order of magnitude estimate empirically extracted from Figure 1. In the right panel of Figure 3 we show the resulting one-dimensional (1D) marginal posterior distribution of the parameter μ obtained from the corresponding run. Finally, we have checked that Figures 1 and 3 are quantitatively similar for all other parameters.

As we mentioned above, the IAT gives an estimate of the number of steps between independent samples (it roughly corresponds to the period of oscillation, measured in number of steps, of the chain in the whole range of the parameter). Therefore, in order to have a true unbiased set of independent samples, one has to “thin” the chain with a step size of roughly $\bar{\tau}$. This greatly decreases the statistics available from the MCMC run. Conceptually there is nothing wrong with having correlated samples, provided they are distributed according to the target distribution, however, even though this would increase the effective available statistics, it would generally affect the estimate of the uncertainties in the Bayesian inference [50, 51]. We defer a careful study of the issue of thinning to a forthcoming publication [52], while here we limit ourselves to describe the procedure we followed to get a rich enough sample.

We have run EMCEE3 for $10^6 + 5 \cdot 10^3$ steps with 10^3 walkers for 11 times. From each run we have discarded a pre-run of $5 \cdot 10^3$ steps, which is a few times $\bar{\tau}$, and thinned the chain with a step size of 10^3 , i.e. roughly $\bar{\tau}$.⁷ Each run then delivered 10^6 roughly

⁷Even though the \hat{R}_c analysis we performed suggests robust convergence after few 10^4 steps, considering the length of the samplers we used (10^6 steps) and the large thinning value (10^3 steps), the difference between discarding a pre-run of $5 \cdot 10^3$ versus a few 10^4 steps is negligible. We have therefore set the burn-in number of steps to $5 \cdot 10^3$ to slightly improve the efficiency of our MCMC generation.

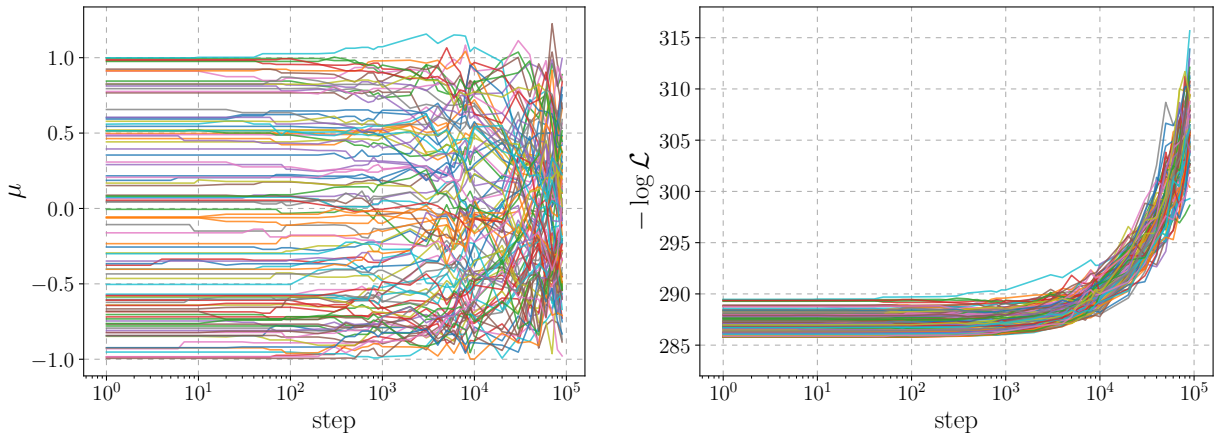


Figure 4: Evolution of the chains in an EMCEE3 sampling of the LF in eq. (11) with 200 walkers and 10^5 steps using the GAUSSIANMOVE algorithm, that updated one parameter at a time, with a variance $5 \cdot 10^{-4}$. The plots show the explored values of the parameter μ (left) and of minus log-likelihood $-\log \mathcal{L}$ (right) versus the number of steps for a random subset of 10^2 of the 200 chains. For visualization purposes, values in the plots are computed only for numbers of steps included in the set $\{a \cdot 10^b\}$ with $a \in [1, 9]$ and $b \in [0, 6]$.

independent samples. With parallelization, the sampler generates and stores about 22 steps per second.⁸ The final sample obtained after all runs consists of $1.1 \cdot 10^7$ samples. We stored 10^6 of them as the test set to evaluate our DNN models, while the remaining 10^7 are used to randomly draw the different training and validation sets used in the following.

2. Biased sample S_2

The second sampling has been used to enrich the training and test sets with points corresponding to large values of the LF, i.e. points close to the maximum for each fixed value of μ . In this case we initialised 200 walkers in maxima of the LF calculated for random values of μ , extracted according to a uniform probability distribution in the interval $[-1, 1]$.⁹ Moreover, the proposals have been updated using a Gaussian random

⁸All samplings presented in the paper were produced with a SYS-7049A-T Supermicro[®] workstation configured as follows: Dual Intel[®] Xeon[®] Gold 6152 CPUs at 2.1GHz (22 physical cores), 128 Gb of 2666 MHz Ram, Dual NVIDIA[®] RTX 2080-Ti GPUs and 1.9Tb M.2 Samsung[®] NVMe PM963 Series SSD (MZ1LW1T9HMLS-00003). Notice that speed, in our case, was almost constant for a wide choice of the number of parallel processes in the range $\sim 30 - 88$, with CPU usage never above about 50%. We therefore conclude that generation speed was, in our case, limited by data transfer and not by CPU resources, making parallelization less than optimally efficient.

⁹This interval has been chosen smaller than the interval of μ considered in the unbiased sampling since values of μ outside this interval correspond to values of the LF much smaller than the global maximum, that are not relevant from the frequentist perspective. The range for the biased sampling can be chosen a posteriori by looking at the frequentist confidence intervals on μ .

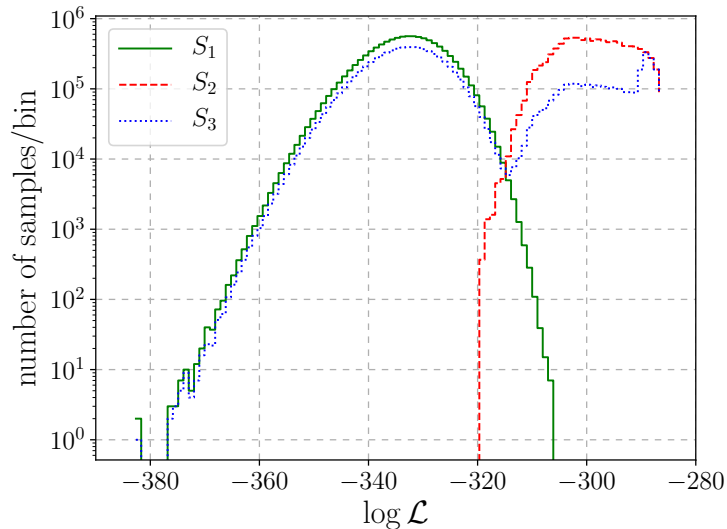


Figure 5: Distribution of $\log \mathcal{L}$ values in S_1 , S_2 and S_3 . S_1 represents the unbiased sampling, S_2 is constructed close to the maximum of $\log \mathcal{L}$, and S_3 is obtained mixing the previous two as explained in the text.

move with variance $5 \cdot 10^{-4}$ (small moves) of a single parameter at a time. In this way, the sampler starts exploring the region of parameters corresponding to local maxima, i.e. large values of the LF, and then slowly moves towards the tails. Once the LF gets further and further from the local maxima, the chains do not explore the region of local maxima anymore. Therefore, in this case we do not want to discard a pre-run, neither to check convergence, which implies that this sampling will have a strong bias (obviously, since we forced the sampler to explore only a particular region).

In Figure 4 we show the evolution of the chains for the parameter μ (left panel) together with the corresponding values of $\log \mathcal{L}$ (right panel) for an illustrative (random) set of 100 chains. Comparing Figure 4 with Figure 1, we see that now the moves of each chain are much smaller and the sampler generates many points all around the local maxima at which the chains are initialised.

In order to ensure a rich enough sampling close to local maxima of the LF, we have made 10^5 iterations for each walker. Since moves are much smaller than in the previous case (only one parameter is updated at a time), the efficiency in this case is very large, $\varepsilon \approx 1$. We could therefore obtain a sampling of $1.1 \cdot 10^7$ points by randomly picking points within the $10^5 \cdot 200 \cdot \varepsilon$ samples. As for S_1 , two samples of 10^6 and 10^7 points have been stored separately: the first serves to build the test set, while the second is used to construct the training and validation sets. As mentioned before, this is a biased sample, and therefore should only be used to enrich the training sample to properly learn the LF close to the maximum (and to check results of the frequentist analysis),

but it cannot be used to make any posterior inference. Due to the large efficiency, this sampling took less than one hour to be generated.

3. Mixed sample S_3

The mixed sample S_3 is built from S_1 and S_2 in order to properly populate both the large probability mass region and the large log-likelihood region. Moreover, we do not want a strong discontinuity for intermediate values of the LF, which could become relevant, for instance, when combining with another analysis that prefers slightly different values of the parameters. For this reason, we have ensured that also intermediate values of the LF are represented, even though with a smaller effective weight, and that no more than a factor of 100 difference in density of examples is present in the whole region $-\log \mathcal{L} \in [285, 350]$. Finally, in order to ensure a good enough statistics close to the maxima, we have enriched further the sample above $\log \mathcal{L} \approx -290$ (covering the region $\Delta \log \mathcal{L} \lesssim 5$).

S_3 has been obtained taking all samples from S_2 with $\log \mathcal{L} > -290$ (around 10% of all samples in S_2), 70% of samples from S_1 (randomly distributed), and the remaining fraction, around 20%, from S_2 with $\log \mathcal{L} < -290$. With this procedure we obtained a total of $10^7(10^6)$ train(test) samples. We have checked that results do not depend strongly on the assumptions made to build S_3 , provided enough examples are present in all the relevant regions in the training sample.

The distribution of the LF values in the three samples are shown in Figure 5 (for the 10^7 points in the training/validation set).

We have used the three samples as follows: examples drawn from S_3 were used to train the full DNNLikelihood, while results have been checked against S_1 in the case of Bayesian posterior estimations and against S_2 (together with results obtained from a numerical maximisation of the analytical LF) in the case of frequentist inference. Moreover, we also present a “Bayesian only” version of the DNNLikelihood, trained using only points from S_1 .

3.1.1 Bayesian inference

In the Bayesian approach one is interested in marginal distributions, used to compute marginal posterior probabilities and credibility intervals. For instance, in the case at hand, one may be interested in two-dimensional (2D) marginal probability distributions in the parameter space $(\mu, \boldsymbol{\delta})$, such as

$$p(\mu, \delta_i) = \int d\delta_1 \cdots \int d\delta_{i-1} \int d\delta_{i+1} \cdots \int d\delta_{95} \mathcal{L}(\mu, \boldsymbol{\delta}) \pi(\mu), \quad (14)$$

or in 1D HPDI corresponding to probabilities $1 - \alpha$, such as

$$1 - \alpha = \int_{\mu_{\text{low}}}^{\mu_{\text{high}}} d\mu \int d\delta_i p(\mu, \delta_i). \quad (15)$$

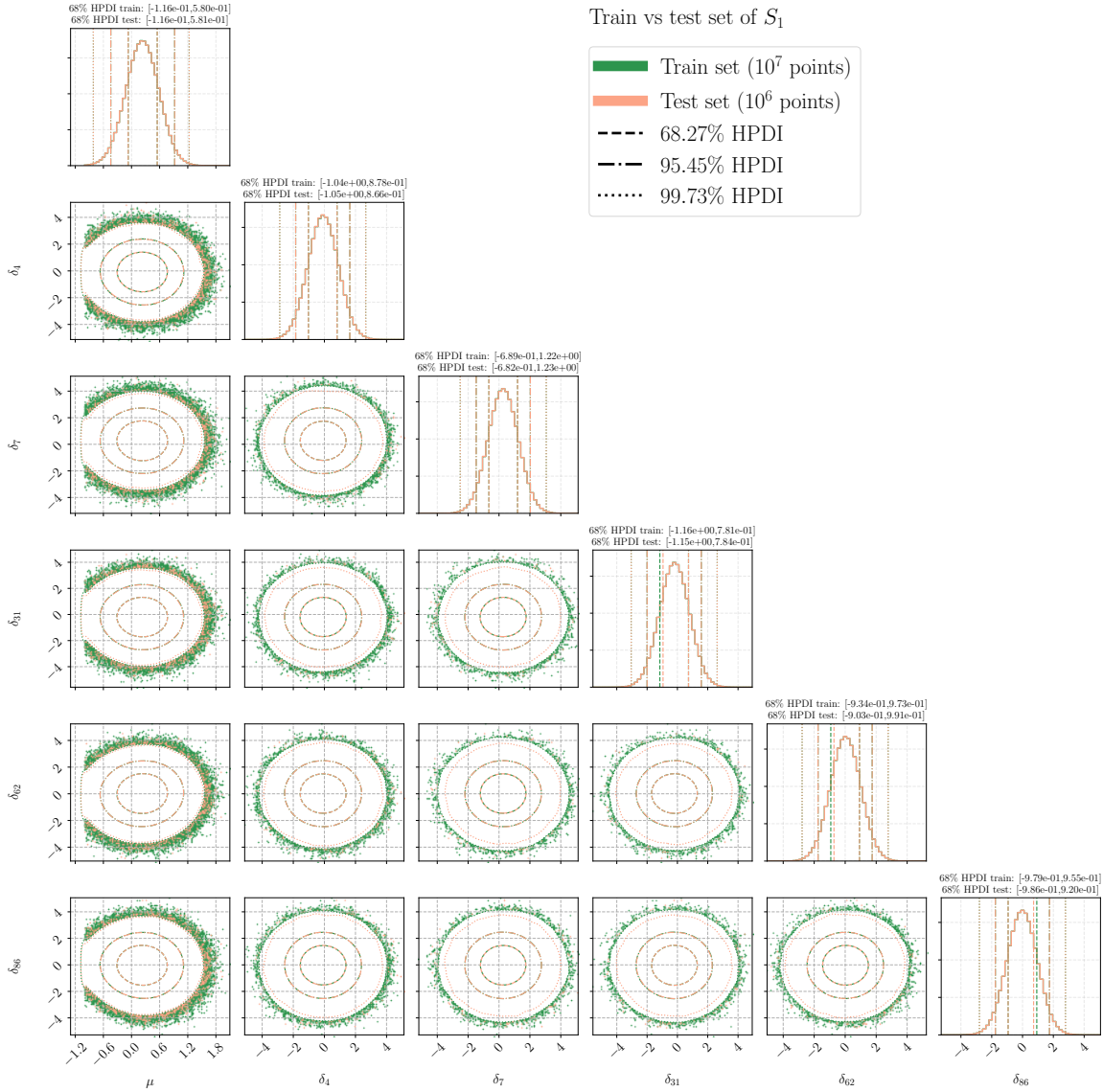


Figure 6: 1D and 2D posterior marginal probability distributions for a subset of parameters from the unbiased S_1 . This gives a graphical representation of the sampling obtained through MCMC. The green (darker) and red (lighter) points and curves correspond to the training set (10^7 points) and test set (10^6 points) of S_1 , respectively. Histograms are made with 50 bins and normalised to unit integral. The dotted, dot-dashed, and dashed lines represent the 68.27%, 95.45%, 99.73% 1D and 2D HPDI. The difference between green (darker) and red (lighter) lines gives an idea of the uncertainty on the HPDI due to finite sampling. Numbers for the 68.27% HPDI for the parameters in the two samples are reported above the 1D plots.

HPDI	$\mu > -1$	$\mu > 0$
68.27%	$[-0.12, 0.58]$	0.48
95.45%	$[-0.47, 0.92]$	0.86
99.73%	$[-0.82, 1.26]$	1.22

Table 1: HPDIs obtained using all 10^7 samples from the training set of S_1 . The result is shown both for $\mu > -1$ and $\mu > 0$ (only the upper bound is given in the latter case).

All these integrals can be discretized and computed by just summing over quantities evaluated on a proper unbiased LF sampling.

This can be efficiently done with MCMC techniques, such as the one described in Section 3.1. For instance, using the sample S_1 we can directly compute HPDIs for the parameters. Figure 6 shows the 1D and 2D posterior marginal probability distributions of the subset of parameters $(\mu, \delta_4, \delta_7, \delta_{31}, \delta_{62}, \delta_{86})$ obtained with the training set (10^7 points, green (darker)) and test set (10^6 points, red (lighter)) of S_1 . Figure 6 also shows the 1D and 2D 68.27%, 95.45%, 99.73% HPDIs. All the HPDIs, including those shown in Figure 6, have been computed by binning the distribution with 50 bins, estimating the interval, and increasing the number of bins by 50 until the interval splits due to statistical fluctuations.

The results for μ with the assumptions $\mu > -1$ and $\mu > 0$, estimated from the training set, which has the largest statistics, are given in Table 1. Figure 6 shows how the 1D marginal probability distributions are extremely accurate up to the 99.73%, while the 2D ones, for the same interval, start showing differences, due to the finite sample size.

Considering that the sample sizes used to train the DNN range from 10^5 to $5 \cdot 10^5$, we do not consider probability intervals higher than 99.73%. Obviously, if one is interested in covering higher HPDIs, larger training sample sizes need to be considered (for instance, to cover a Gaussian 5σ interval, that corresponds to a probability $1 - 5.7 \cdot 10^{-7}$, even only on the 1D marginal distributions, a sample with $\gg 10^7$ points would be necessary). We will not consider this case in the present paper.

3.1.2 Frequentist inference

In a frequentist inference one usually constructs a test statistics $\lambda(\mu, \boldsymbol{\theta})$ based on the LF ratio

$$\lambda(\mu, \boldsymbol{\delta}) = \frac{\mathcal{L}(\mu, \boldsymbol{\delta})}{\mathcal{L}_{\max}(\hat{\mu}, \hat{\boldsymbol{\delta}})}. \quad (16)$$

Since one would like the test statistics to be independent of the nuisance parameters, it is common to use instead the profiled likelihood, obtained replacing the LF at each value of μ with its maximum value (over the nuisance parameters volume) for that value of μ . One

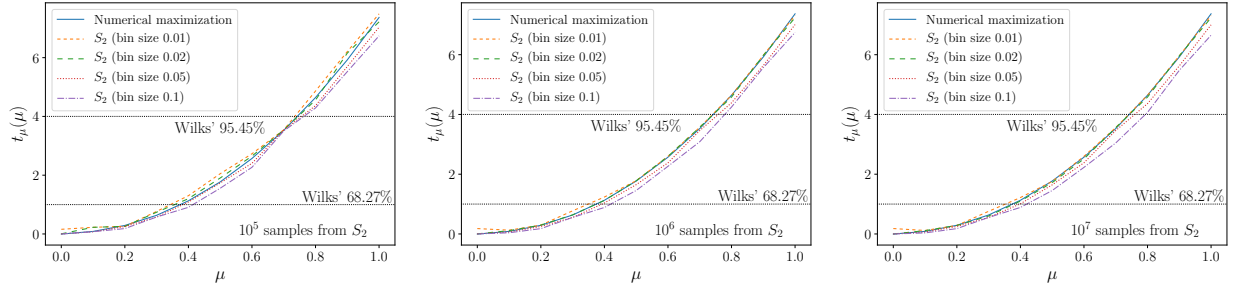


Figure 7: Comparison of the t_μ test-statistics computed using numerical maximisation of eq. (17) and using a variable sample size from S_2 . We show the result obtained searching from the maximum by using different binning in t_μ with bin size 0.01, 0.02, 0.05, 0.1 around each value of μ (between 0 and 1 in steps of 0.1).

can then construct a test statistics t_μ based on the profiled (log)-likelihood ratio, given by

$$t_\mu = -2 \log \frac{\mathcal{L}_{\text{prof}}(\mu)}{\mathcal{L}_{\text{max}}} = -2 \log \frac{\sup_{\delta} \mathcal{L}(\mu, \delta)}{\sup_{\mu, \delta} \mathcal{L}(\mu, \delta)} = -2 \left(\sup_{\delta} \log \mathcal{L}(\mu, \delta) - \sup_{\mu, \delta} \log \mathcal{L}(\mu, \delta) \right). \quad (17)$$

Whenever suitable general conditions are satisfied, and in the limit of large data sample, by Wilks' theorem [53] the distribution of this test-statistics approaches a χ^2 distribution that is independent of the nuisance parameters δ and has a number of degrees of freedom equal to $\dim \mathcal{L} - \dim \mathcal{L}_{\text{prof}}$ [54]. In our case t_μ can be computed using numerical maximisation on the analytic LF, but it can also be computed from S_2 (and S_3 , which is identical in the large likelihood region), which was constructed with the purpose of describing the LF as precisely as possible close to local maxima. In Figure 7 we show the result for t_μ using both approaches for different sample sizes drawn from S_2 . The three samples from S_2 used for the maximisation, with sizes 10^5 , 10^6 , and 10^7 (full training set of S_2), contain in the region $\mu \in [0, 1]$ around $5 \cdot 10^4$, $5 \cdot 10^5$, and $5 \cdot 10^6$ points respectively, which results in increasing statistics in each bin and a more precise and stable prediction for t_μ . As it can be seen 10^5 points, about half of which contained in the range $\mu \in [0, 1]$, are already sufficient, with a small bin size of 0.02, to reproduce the t_μ curve with great accuracy. As expected, larger bin sizes result in too high local maxima estimates, leading to an underestimate of t_μ .

Under Wilks' theorem assumptions, t_μ should be distributed as a χ^2_1 (1 d.o.f.) distribution, from which we can determine CL upper limits. The 68.27%(95.45%) CL upper limit (under the Wilks' hypotheses) is given by $t_\mu = 1(4)$, corresponding to $\mu < 0.37(0.74)$. These upper limits are compatible with the ones found in ref. [9], and are quite smaller than the corresponding upper limits of the HPDI obtained with the Bayesian analysis in Section 3.1.1 (see Table 1). This may suggest that the result obtained using the asymptotic approximation for t_μ is underestimating the upper limit (undercoverage). This is somehow expected for the search under consideration, given the large number of bins in which the observed number of events is below 5 or even 3 (see Figure 2 of ref. [9]). Indeed, the true distribution

of t_μ is expected to depart from a χ_1^2 distribution when the hypotheses of Wilks’ theorem are violated. The study of the distribution of t_μ is related to the problem of coverage of frequentist confidence intervals, and requires to perform pseudo-experiments. We present results on the distribution of t_μ obtained through pseudo-experiments in Appendix B. The important conclusion is that using the distribution of t_μ generated with pseudo-experiments, CL upper limits become more conservative by up to almost a factor of two, depending on the choice of the approach used to treat nuisance parameters. This shows that the upper limits computed through asymptotic statistics undercover, in this case, the actual upper bounds on μ .

4 The DNNLikelihood

The sampling of the full likelihood discussed above has been used to train a DNN regressor constructed from multiple fully connected layers, i.e. a multilayer perceptron (MLP). The regressor has been trained to predict values of the LF given a vector of inputs made by the physical and nuisance parameters. In order to introduce the main ingredients of our regression procedure and DNN training, we first show how models trained using only points from S_1 give reliable and robust results in the case of the Bayesian approach. Then we discuss the issue of training with samples from S_3 to allow for maximum likelihood based inference. Finally, once a satisfactory final model is obtained, we show again its performance for posterior Bayesian estimates.

4.1 Model architecture and optimisation

We used KERAS [55] with TENSORFLOW [56] backend, through their PYTHON implementation, to train a MLP and considered the following hyperparameters to be optimised, the value of which defines what we call a *model* or a DNNLikelihood.

- **Size of training sample**

In order to assess the performance of the DNNLikelihood given the training set size we considered three different values: 10^5 , $2 \cdot 10^5$ and $5 \cdot 10^5$. The training set (together with an half sized evaluation set) has been randomly drawn from S_1 for each model training, which ensures the absence of correlation between the models due to the training data: thanks to the large size of S_1 (10^7 samples) all the training sets can be considered roughly independent. In order to allow for a consistent comparison, all models trained with the same amount of training data have been tested with a sample from the test set of S_1 , and with half the size of the training set. In general, and in particular in our interpolation problem, increasing the size of the training set allows to reduce the generalization error and therefore to obtain the desired performance on the test set.

- **Loss function**

In Section 2 we have argued that both MAE and MSE are suitable loss functions to

learn the log-likelihood function. In our optimisation procedure we tried both, always finding (slightly) better results for the MSE. We therefore choose the MSE as our loss function in all results presented here.

- **Number of hidden layers**

From a preliminary optimisation we concluded that more than a single HL (deep network) always performs better than a single HL (shallow network). However, in the case under consideration, deeper networks do not seem to perform much better than 2HL networks, even though they are typically much slower to train and to make predictions. Therefore, after this preliminary assessment, we focused on 2HL architectures.

- **Number of nodes on hidden layers**

We considered architectures with the same number of nodes on the two hidden layers. The number of trainable parameters (weights) in the case of n fully connected HLs with the same number of nodes d_{HL} is given by

$$d_{\text{HL}} (d_{\text{input}} + (n - 1)d_{\text{HL}} + (n + 1)) + 1, \quad (18)$$

where d_{input} is the dimension of the input layer, i.e. the number of independent variables, 95 in our case. DNNs trained with stochastic gradient methods tend to small generalization errors even when the number of parameters is larger than the training sample size [57]. Overfitting is not an issue in our interpolation problem [58]. In our case we considered HLs not smaller than 500 nodes, which should ensure enough *bandwidth* throughout the network and model *capacity*. In particular we compared results obtained with 500, 1000, 2000, and 5000 nodes on each HL, corresponding to 299001, 1098001, 4196001, and 25490001 trainable parameters.

- **Activation function on hidden layers**

We compared RELU [59], ELU [60], and SELU [61] activation functions and the latter one showed to fit better our problem. In order to correctly implement the SELU activation in Keras we initialised all weights using the KERAS “lecun_normal” initialiser [61, 62].

- **Batch size**

When using a stochastic gradient optimisation technique, of which ADAM is an example, the minibatch size is an hyperparameter. For the training to be stochastic, the batch size should be much smaller than the training set size, so that each minibatch can be considered roughly independent. Large batch sizes lead to more accurate weight updates and, due to the parallel capabilities of GPUs, to faster training time. However, smaller batch sizes usually contribute to regularize and avoid overfitting. After a preliminary optimisation obtained changing the batch size from 256 to 4096, we concluded that the best performances were obtained by keeping the number of batches roughly fixed to 200 when changing the training set size. In particular, choosing batch sizes among powers of two, we have used 512, 1024 and 2048 for 10^5 , $2 \cdot 10^5$ and $5 \cdot 10^5$

training set sizes respectively. Notice that increasing the batch size when enlarging the training set, also allowed us to keep the initial learning rate fixed [63]. Similar results could be obtained by keeping a fixed batch size of 512 and reducing the starting learning rate when enlarging the training set.

- **Optimiser**

We used the ADAM optimiser with default parameters, and in particular with learning rate $\epsilon = 0.001$. We reduced the learning rate by a factor 0.2 every 40 epochs without improvements on the validation loss within an absolute amount (MIN_DELTA in KERAS) $1/N_{\text{points}}$, with N_{points} the training set size. Indeed, being the KERAS MIN_DELTA parameter absolute and not relative to the value of the loss function, we needed to reduce it when getting smaller losses (better models). We have found that $1/N_{\text{points}}$ corresponded roughly to one to few permil of the best minimum validation loss obtained for all different training set sizes. This value turned out to give the best results with reasonably low number of epochs (fast enough training). Finally, we performed early stopping [64, 65] using the same MIN_DELTA parameter and no improvement in the validation loss for 50 epochs. This ensured that training did not go on for too long without substantially improving the result. We also tested the newly proposed ADABOUND optimiser [66] without seeing, in our case, large differences.

Notice that the process of choosing and optimising a model depends on the LF under consideration (dimensions, number of modes, etc.) and this procedure should be repeated for different LFs. However, good initial points for the optimisation could be chosen using experience from previously constructed DNNLikelihoods.

As we discussed in Section 2, there are several metrics that we can use to evaluate our model. Based on the results obtained by re-sampling the DNNLikelihood with EMCEE3, we see a strong correlation between the quality of the re-sampled probability distribution (i.e. of the final Bayesian inference results) and the metric corresponding to the median of the K-S test on the 1D posterior marginal distributions. We therefore present results focusing on this evaluation metric. When dealing with the Full DNNLikelihood trained with the biased sampling S_3 we also consider the performance on the mean relative error on the predicted t_μ test statistics when choosing the best models.

4.2 The Bayesian DNNLikelihood

From a Bayesian perspective, the aim of the DNNLikelihood is to be able, through a DNN interpolation of the full LF, to generate a sampling analog to S_1 , which allows to produce Bayesian posterior density distributions as close as possible to the ones obtained using the true LF, i.e. the S_1 sampling. Moreover, independently on how complicated to evaluate the original LF is, the DNNLikelihood is extremely fast to compute, allowing for very fast sampling.¹⁰ The EMCEE3 MCMC package allows, through vectorization of the input function

¹⁰In this case the original likelihood is extremely fast to evaluate either, since it is known in analytical form. This is usually not the case in actual experimental searches involving theory and detector simulations.

for the log-probability, to profit of parallel GPU predictions, which made sampling of the DNNLikelihood roughly as fast as the original analytic LF.

We start by considering training using samples drawn from the unbiased S_1 . The independent variables all vary in a reasonably small interval around zero and do not need any preprocessing. However, the $\log \mathcal{L}$ values in S_1 span a range between around -380 and -285 . This is both pretty large and far from zero for the training to be optimal. For this reason we have pre-processed data scaling them to zero mean and unit variance. Obviously, when predicting values of $\log \mathcal{L}$ we applied the inverse function to the DNN output.

We rank models trained during our optimisation procedure by the median p -value of 1D K-S test on all coordinates between the test set and the prediction performed on the validation set. The best models are those with the highest median p -value. In Table 2 we show results for the best model we obtained for each training sample size. Results have been obtained by training 5 identical models and taking the best one. We call these four best models $B_1 - B_3$ (B stands for Bayesian). All three models have two HLs with $5 \cdot 10^3$ nodes each, and are therefore the largest we consider in terms of number of parameters. However, it should be clear that the gap with smaller models is extremely small in some cases with some of the models with less parameters in the ensemble of 5 performing better than some others with more parameters. This also suggests that results are not too sensitive to model dimension, making the DNNLikelihood pretty robust.

Name	B_1	B_2	B_3
Sample size ($\times 10^5$)	1	2	5
Epochs	178	268	363
Loss train (MSE) ($\times 10^{-3}$)	0.14	0.088	0.054
Loss val (MSE) ($\times 10^{-3}$)	10.11	6.66	3.9
Loss test (MSE) ($\times 10^{-3}$)	10.02	6.64	3.9
ME train ($\times 10^{-3}$)	0.47	0.53	0.28
ME val ($\times 10^{-3}$)	5.44	2.58	1.76
ME test ($\times 10^{-3}$)	4.91	2.31	1.72
Median p -value of 1D K-S test vs pred. on train	0.41	0.46	0.39
Median p -value of 1D K-S test vs. pred. on val.	0.24	0.33	0.43
Median p -value of 1D K-S val vs. pred. on test	0.24	0.40	0.34
Training time (s)	1007	2341	8446
Prediction time (μ s/point)	11.5	10.4	14.5

Table 2: Results for the best models (Bayesian DNNLikelihood) for different training sample size. All models have been trained for 5 times to check the stability of the result and the best performing one has been quoted. Prediction time is evaluated on a test set with half the size of the training set using the same batch size used in training, and evaluating on a Nvidia Tesla V100 GPU with 32GB of RAM. All best models have $d_{\text{HL}} = 5 \cdot 10^3$.

Figure 8 shows the learning curves obtained for the values of the hyperparameters shown

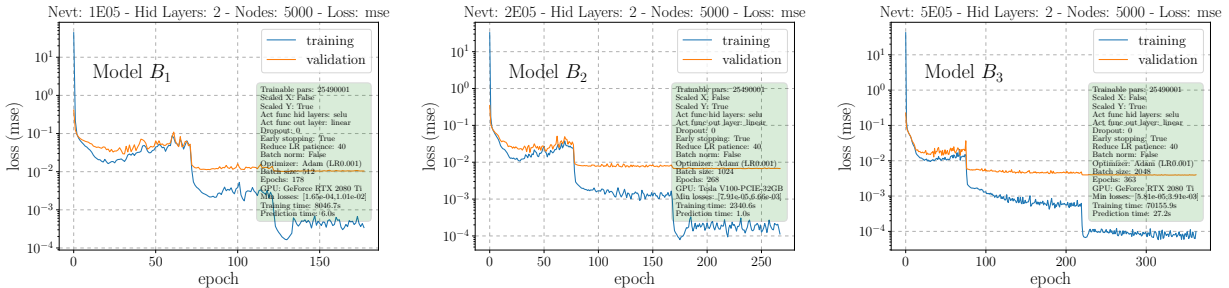


Figure 8: Training and validation loss (MSE) vs number of training epochs for models $B_1 - B_3$. The jumps correspond to points of reduction of the ADAM optimiser learning rate.

in the legends. Early stopping is usually triggered after a few hundred epochs (ranging from around 200 to 500, with the best models around 200 – 300) and values of the validation loss (MSE) that range in the interval $\approx [0.01, 0.003]$. Values of the validation ME, which, as explained in Section 2 correspond to the K-L divergence for the LF, range in the $\approx [1, 5] \cdot 10^{-3}$, which, together with median of the p -value of the 1D K-S tests in the range 0.2 – 0.4 deliver very accurate models. Training times are not prohibitive, and range from less than one hour to a few hours for the models we considered on a Nvidia Tesla V100 GPU with 32GB of RAM. Prediction times, using the same batch sizes used during training, are in the ballpark of 10 – 15 μ s/point, allowing for very fast sampling and inference using the DNNLikelihood. Finally, as shown in Table 2, all models present very good generalization when going from the evaluation to the test set, with the generalization error decreasing with the sample size as expected.

In order to get a full quantitative assessment of the performances of the Bayesian DNN-Likelihood, we compared the results of a Bayesian analysis performed using the test set of S_1 and each of the models $B_1 - B_3$. This was done in two ways. Since the model is usually a very good fit of the LF, we reweighted each point in S_1 using the ratio between the original likelihood and the DNNLikelihood (reweighting). This procedure is so fast that can be done for each trained model during the optimisation procedure giving better insights on the choice of hyperparameters. Once the best model has been chosen, the result of reweighting has been checked by directly sampling the DNNLikelihoods with EMCEE3.¹¹ We present results obtained by sampling the DNNLikelihoods in the form of 1D and 2D marginal posterior density plots on a chosen set of parameters ($\mu, \delta_{10}, \delta_{40}, \delta_{70}, \delta_{95}$).

We have sampled the LF using the DNNLikelihoods $B_1 - B_3$ with the same procedure used for S_1 .¹² Starting from model B_1 (Figure 9), Bayesian inference is well reproduced up

¹¹Sampling has been done on the same hardware configuration mentioned in Footnote 8. However, in this case log-probabilities have been computed in parallel on GPUs (using the "vectorize" option of EMCEE3).

¹²Since the effect of autocorrelation of walkers is much smaller than the intrinsic error of the DNN, to speed up sampling we have used 1024 walkers with 10^5 steps, discarded a burn-in phase of $5 \cdot 10^4$ steps and thinned the remaining $5 \cdot 10^4$ with a step size of 50 to end up with 10^6 samples from each of the DNNLikelihoods.

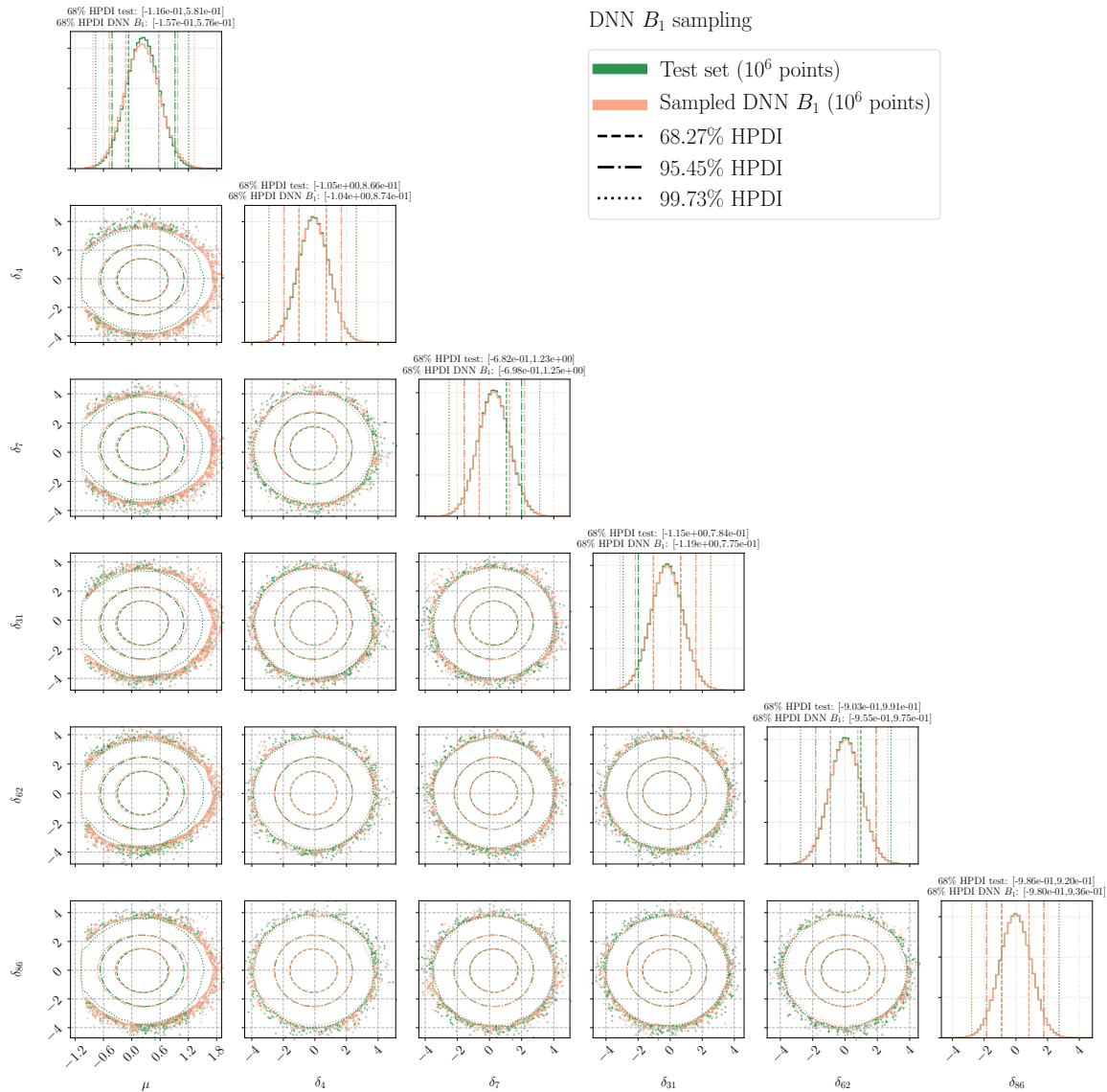


Figure 9: 1D and 2D posterior marginal probability distributions for a subset of parameters from the unbiased S_1 . The green (darker) distributions represent the test set of S_1 , while the red (lighter) distributions are obtained by sampling the DNNLikelihood B_1 . Histograms are made with 50 bins and normalised to unit integral. The dotted, dot-dashed, and dashed lines represent the 68.27%, 95.45%, 99.73% 1D and 2D HPDI. The difference between green (darker) and red (lighter) lines gives an idea of the uncertainty on the HPDI due to finite sampling. Numbers for the 68.27% HPDI for the parameters in the two samples are reported above the 1D plots.

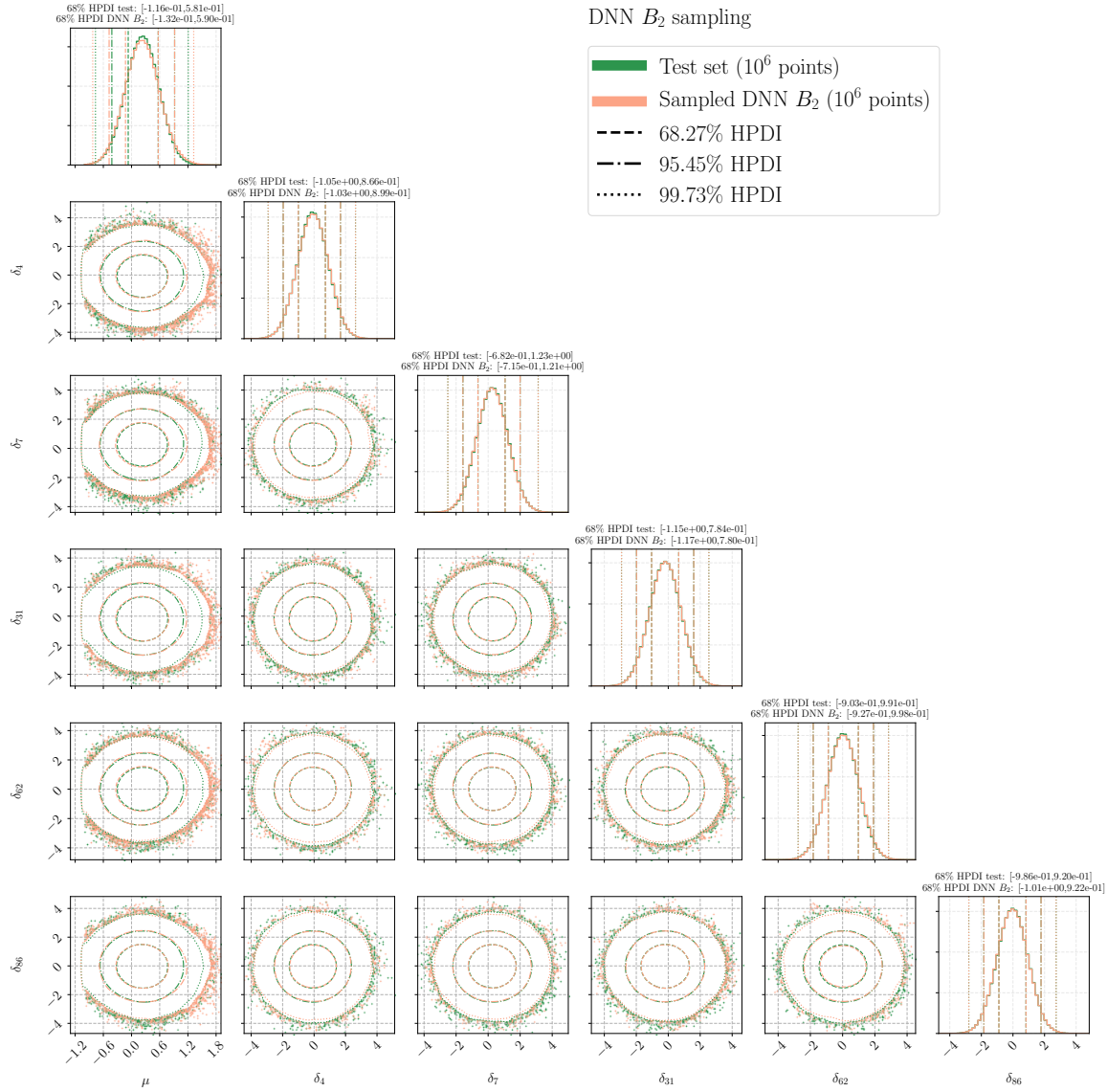


Figure 10: Same as Figure 9 but for the DNNLikelihood B_2 .

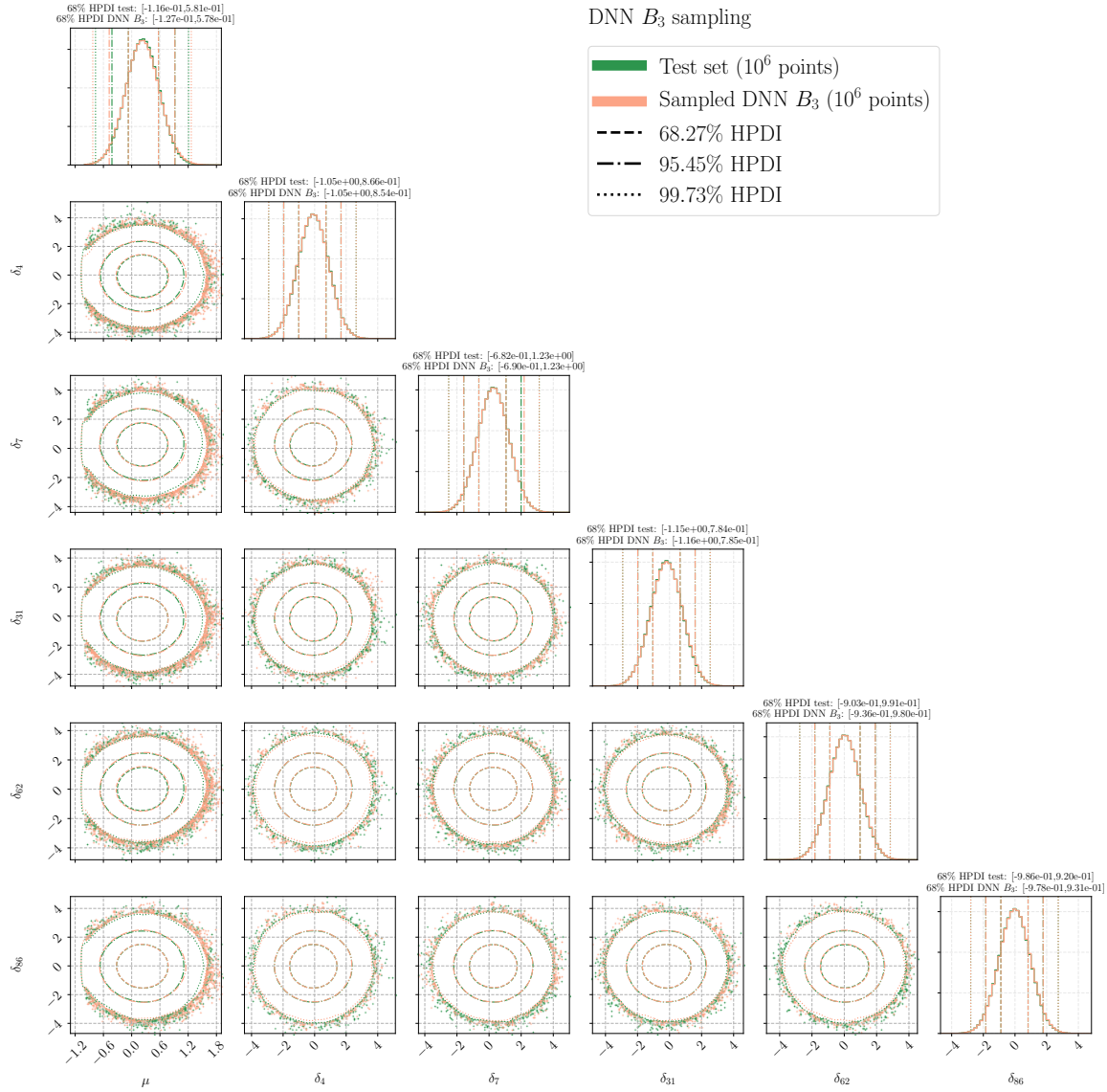


Figure 11: Same as Figure 9 but for the DNNLikelihood B_3 .

HPDI	B_1	B_2	B_3
68.27%	0.50	0.49	0.48
95.45%	0.92	0.91	0.88
99.73%	1.36	1.35	1.29
68.27%	0.49	0.49	0.49
95.45%	0.92	0.91	0.88
99.73%	1.35	1.34	1.29

Table 3: HPDI obtained from the different DNNLikelihood models $B_1 - B_3$ both by reweighting on the test set of S_1 (upper block) and by re-sampling (lower block). Results are only shown as upper bound on μ in the hypothesis $\mu > 0$.

to probability intervals of 95.45%, while large deviations start to arise at 99.73%. This is reasonable, since model B_1 has been trained with only 10^5 points, which are not enough to carefully interpolate in the tails, so that this region is described by the DNNLikelihood through extrapolation. Small deviations also arise for smaller probability intervals, which are again the outcome of a model trained with too few points. Nevertheless, we want to stress that, considering the very small training size and the large dimensionality of the LF, model B_1 already works surprisingly well. This is a common feature of the DNNLikelihood, which, as anticipated, works extremely well in predicting posterior probabilities without the need of a too large training sample, nor a hard tuning of the DNN hyperparameters. When going to models B_2 and B_3 (Figures 10 and 11) predictions become more and more reliable, improving as expected with the number of training points. Notice that the observed deviations are of the same size, and often smaller, than the ones observed in Figure 6 between the larger training set and the smaller test set. Therefore, at least part of the small deviations observed in the DNNLikelihood prediction have to be attributed to the finite size of the training set, and are expected to disappear when further increasing the number of points. Considering the relatively small training and prediction times shown in Table 2, it should be possible, once the desired level of accuracy has been chosen, to enlarge the training and test sets enough to match that precision. For the purpose of this work, we consider the results obtained with models $B_1 - B_3$ already satisfactory, and do not go beyond $5 \cdot 10^5$ training samples.

In order to allow for a fully quantitative comparison, in Table 3 we summarize the Bayesian 1D HPDI obtained with the DNNLikelihoods $B_1 - B_3$ for the parameter μ both using reweighting and re-sampling (only upper bounds for the hypothesis $\mu > 0$). We find that, taking into account the uncertainty arising from our algorithm to compute HPDI (finite binning) and from statistical fluctuations in the tails of the distributions for large probability intervals, the results of Table 3 are in good agreement with those in Table 1. This shows that the Bayesian DNNLikelihood is accurate even with a rather small training sample size of 10^5 and its accuracy quickly improves by increasing the training sample size.

Name	F_1	F_2	F_3
Sample size ($\times 10^5$)	1	2	5
Epochs	183	243	362
Loss train (MSE) ($\times 10^{-3}$)	0.092	0.026	0.030
Loss val (MSE) ($\times 10^{-3}$)	1.18	0.80	0.71
Loss test (MSE) ($\times 10^{-3}$)	1.17	0.80	0.72
ME train ($\times 10^{-3}$)	3.07	0.47	1.1
ME val ($\times 10^{-3}$)	1.78	0.87	0.82
ME test ($\times 10^{-3}$)	1.50	0.68	0.86
Median p -value of 1D K-S test/pred-train	0.53	0.48	0.44
Median p -value of 1D K-S test/pred-val	0.15	0.27	0.20
Median p -value of 1D K-S val/pred-test	0.13	0.31	0.33
Mean error on t_μ	0.11	0.12	0.032
Training time (s)	1236	2819	7114
Prediction time (μ s/point)	11.1	10.8	10.5

Table 4: Results for the best models (Full DNNLikelihood) for different training sample size. All models have been trained for 5 times to check the stability of the result and the best performing has been quoted. Prediction time is evaluated on a test set with half the size of the training set using the same batch size used in training, and evaluating on a Nvidia Tesla V100 GPU with 32GB of RAM. All best models have $d_{\text{HL}} = 5 \cdot 10^3$.

4.3 Frequentist extension and the full DNNLikelihood

We have trained the same model architectures considered for the Bayesian DNNLikelihood using the S_3 sample. In Table 4 we show results for the best models we obtained for each training sample size. Results have been obtained by training 5 identical models and taking the best one. We call these models $F_1 - F_3$ (F may stand for both Full and Frequentist, bearing in mind that these models also allow for Bayesian inference). As we anticipated in the previous Section, the performance gap between models with different number of parameters is very small and often models with less parameters overperform, at least in the hyperparameter space we considered, models with more parameters. This is especially due to our choice of leaving the initial learning rate constant for all architectures, which resulted in a slightly too large learning rate for the bigger models, with a consequently less stable training phase.

This can be seen in Figure 12, where we show the learning curves obtained for the values of the hyperparameters shown in the legends. The training curves for models F_1 and F_2 are less regular than those of model F_3 . Nevertheless, as the learning rate gets reduced, they quickly reach a good validation loss, which promoted them best models for 10^5 and $2 \cdot 10^5$ training set sizes. This did not happen for models trained with $5 \cdot 10^5$ points, which could have benefited from reducing further the initial learning rate. However, we stress that no strong fine-tuning is needed for the DNNLikelihood to perform extremely well, so that we have chosen the best model F_3 with less parameters without pushing optimisation further.

Final results are generally very similar to the ones obtained for the Bayesian DNNLikelihood, with differences arising from the new region of LF values that is learnt from the DNN.

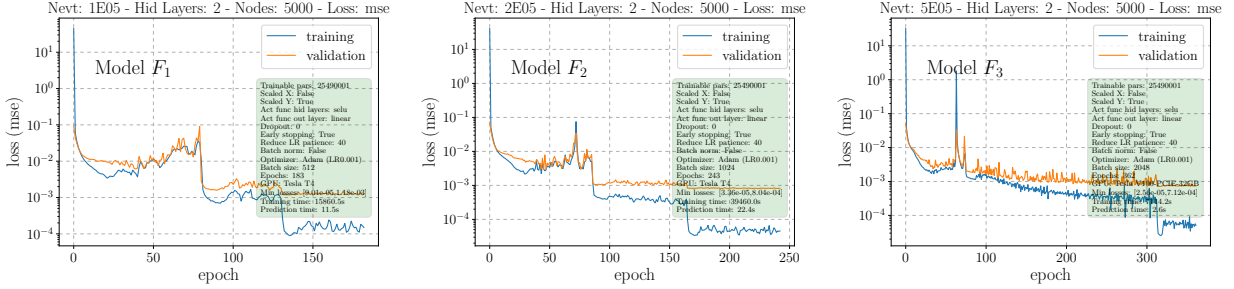


Figure 12: Training and validation loss (MSE) vs number of training epochs for models $F_1 - F_3$. The jumps correspond to points of reduction of the ADAM optimiser learning rate.

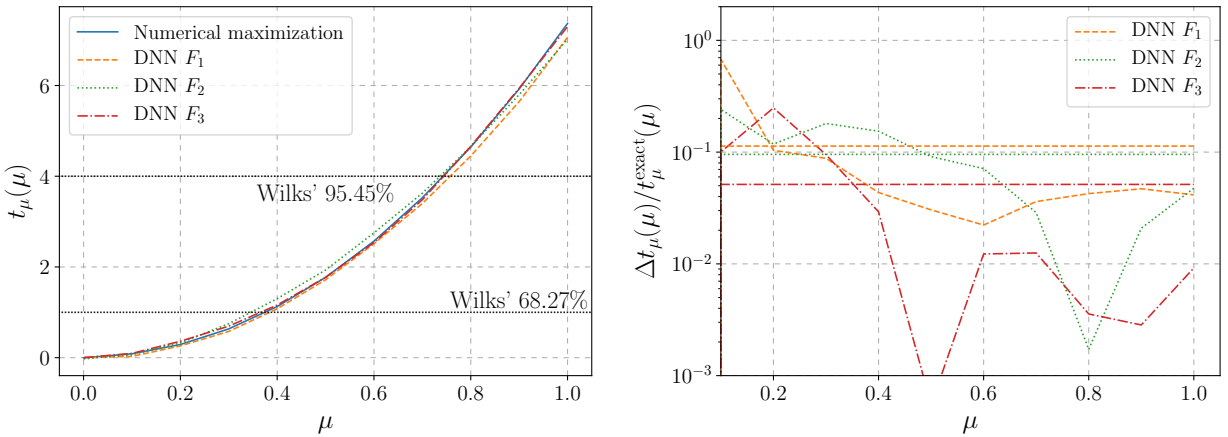


Figure 13: Comparison of the t_μ test-statistics computed using numerical maximisation of the analytic likelihood and of the DNNLikelihoods $F_1 - F_3$. Each of the t_μ prediction from the DNNLikelihoods corresponds to the best out of five trained models. The left and right plots show respectively the t_μ test-statistics and the percentage error computed as $|t_\mu^{\text{DNN}} - t_\mu^{\text{exact}}|/t_\mu^{\text{exact}}$. Horizontal lines in the right plot represent the mean relative error.

Before repeating the Bayesian analysis presented for the Bayesian DNNLikelihood, we present results of frequentist inference using the Full DNNLikelihood. We used the models $F_1 - F_3$ to evaluate the test statistics t_μ . The left panel of Figure 13 shows t_μ obtained using the Full DNNLikelihoods, while the right panel of the same Figure shows the relative error with respect to numerical maximisation of the analytic LF. Apart for some visible deviation in the prediction of model F_1 , it is clear that the Full DNNLikelihood is perfectly able to reproduce the test statistics, allowing for robust frequentist inference, already starting

with relatively small training sets of few hundred thousand points. Clearly, the larger the training set, the smaller the error, with an average percentage error on t_μ (dashed lines in the right panel of Figure 13) that gets as low as a few percent for our models. We found that ensemble learning can help in reducing differences further (keeping under control statistical fluctuations in the training set and in the DNN weights). However, in the particular case under consideration, results are already satisfactory by taking the best out of five identical models trained with random subsets of the training set. For this reason we do not expand on ensemble learning in this work and just consider it as a tool to improve performance in cases where the LF is extremely complicated or has very high dimensionality.

We now show how the full DNNLikelihood is also able to catch all the features necessary for Bayesian inference. We do so by repeating the analysis done for models $B_1 - B_3$ for the full DNNLikelihoods $F_1 - F_3$. Plots obtained by sampling the DNNLikelihood are shown in Figs. 14-16. Results are quantitatively unchanged with respect to the Bayesian DNNLikelihood. For completeness we also summarize in Table 5 the Bayesian 1D HPDI obtained with the DNNLikelihoods $F_1 - F_3$ for the physics parameter μ . Results are compatible with what we found in Table 3 for the Bayesian DNNLikelihood.

HPDI	F_1	F_2	F_3
68.27%	0.50	0.50	0.49
95.45%	0.93	0.92	0.89
99.73%	1.39	1.37	1.31
68.27%	0.50	0.50	0.48
95.45%	0.93	0.93	0.88
99.73%	1.35	1.37	1.28

Table 5: HPDI obtained from the different DNNLikelihood models $F_1 - F_3$ both by reweighting on the test set of S_1 (upper block) and by re-sampling (lower block). Results are only shown as upper bound on μ in the hypothesis $\mu > 0$.

5 Conclusion and future work

Publishing and distributing likelihoods in a simple yet general way is becoming a key issue in many fields, and in particular in the physics of fundamental interactions, where experimental (and phenomenological) results typically involve complicated (and often multi-modal or degenerate) likelihoods which depend on hundreds of parameters. We have introduced the DNNLikelihood framework, in which the full likelihood, binned or unbinned, including the dependence on all nuisance parameters, is published by providing a suitably trained DNN predictor. Distributing the results of experimental or phenomenological analyses using the DNNLikelihood framework allows for the combination of different analyses, for the reinterpretation of the results under different hypotheses, and for the use of the likelihood in a different statistical framework, without loss of information.

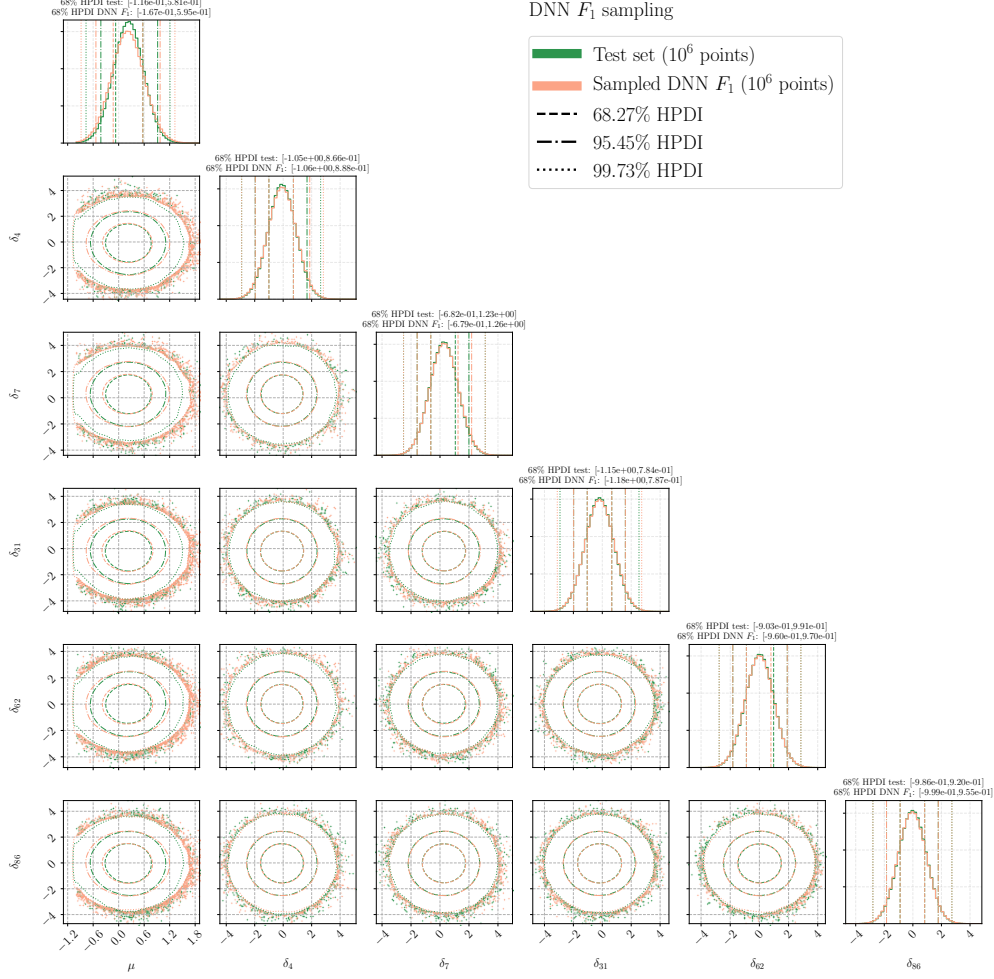


Figure 14: 1D and 2D posterior marginal probability distributions for a subset of parameters from the unbiased S_1 . The green (darker) distributions represent the test set of S_1 , while the red distributions are obtained by sampling the DNNLikelihood F_1 . Histograms are made with 50 bins and normalised to unit integral. The dotted, dot-dashed, and dashed lines represent the 68.27%, 95.45%, 99.73% 1D and 2D HPDI. The difference between green (darker) and red (lighter) lines gives an idea of the uncertainty on the HPDI due to finite sampling. Numbers for the 68.27% HPDI for the parameters in the two samples are reported above the 1D plots.

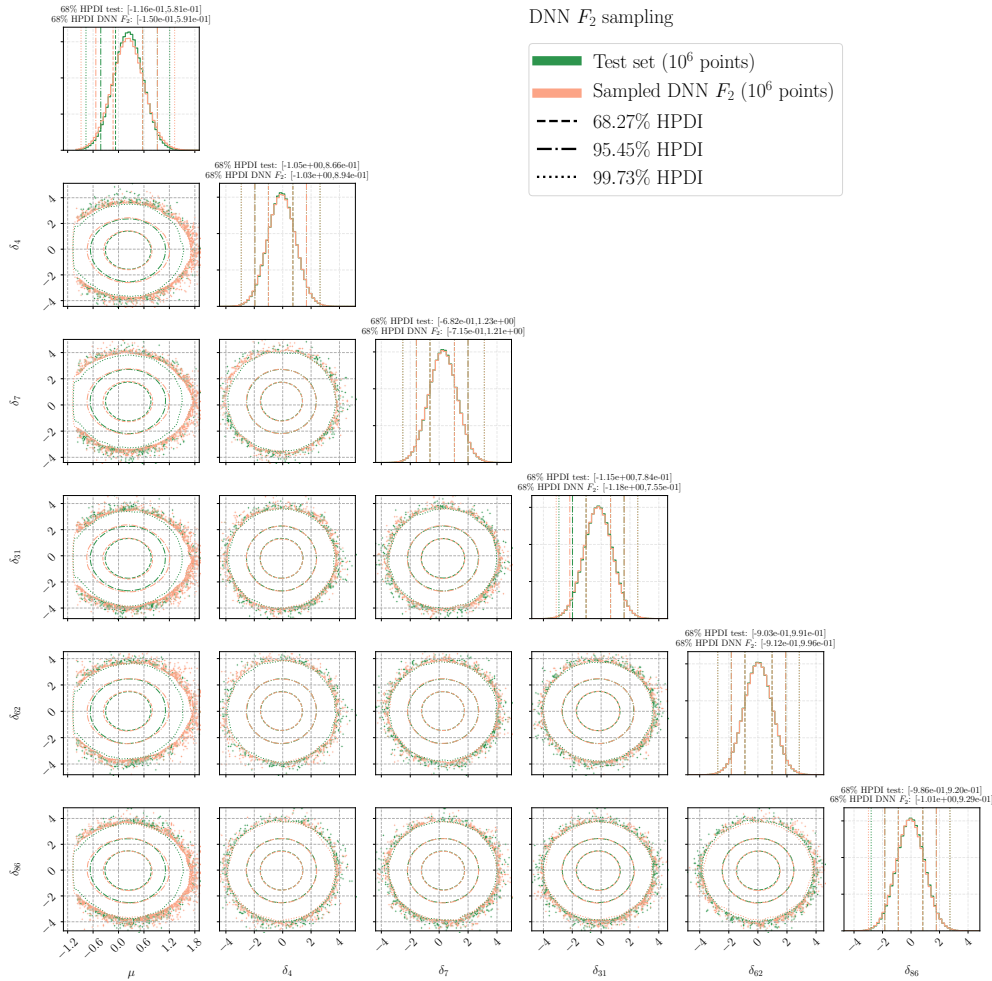


Figure 15: Same as Figure 14 but for the DNNLikelihood F_2 .

We have illustrated the power of the DNNLikelihood discussing in detail the toy experiment presented in ref. [9], which mimics a realistic LHC-like NP search. We found that the DNNLikelihood is able to catch the main features of the true LF, allowing for both frequentist and Bayesian inference, already with a limited amount of training data.

A JUPYTER notebook, together with PYTHON source files which allow to reproduce all results presented above are available on GitHub [\[9\]](#). A dedicated PYTHON package allowing to sample LFs and to build, optimize, train, and store the corresponding DNNLikelihoods is in preparation. This will allow not only allow to construct and distribute DNNLikelihoods, but also to use them for inference both in the Bayesian and frequentist frameworks.

We plan to release soon the first examples of the use of the DNNLikelihood for true experimental likelihoods, including unbinned, multimodal likelihoods, in forthcoming publications.

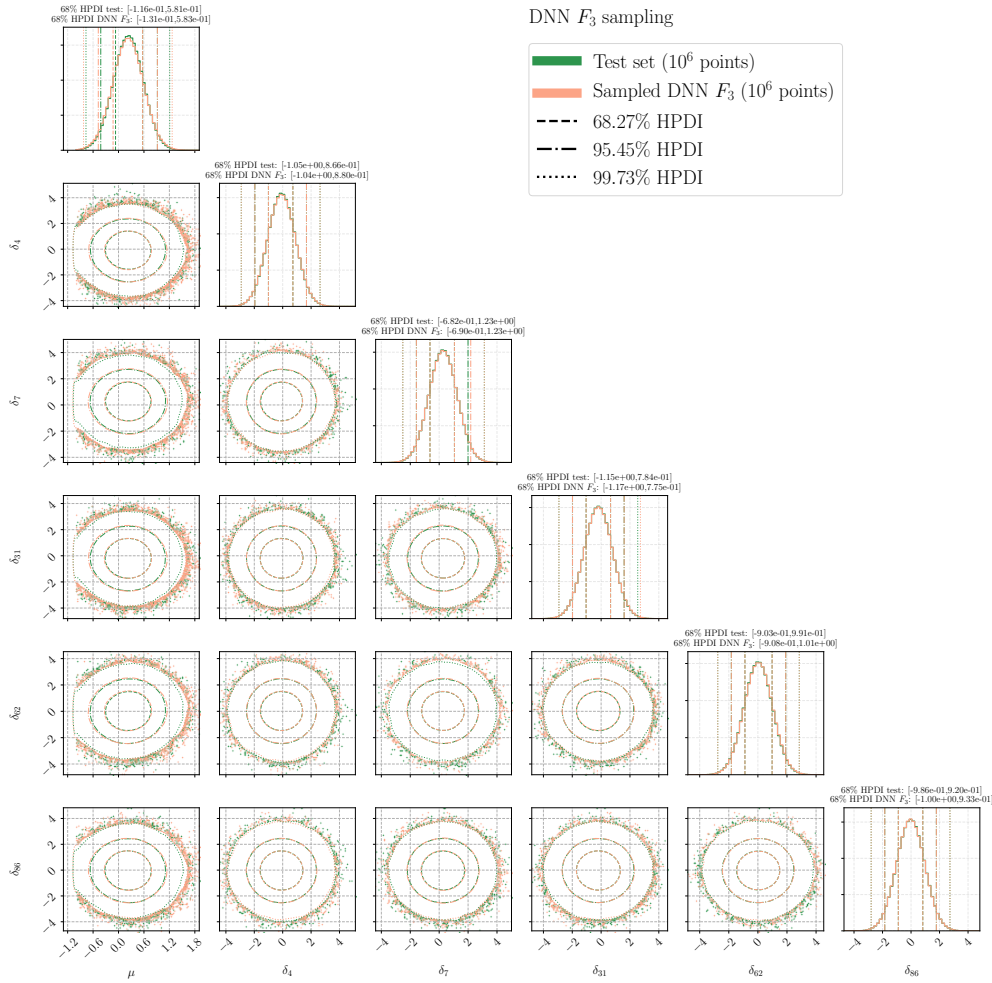


Figure 16: Same as Figure 14 but for the DNNLikelihood F_3 .

As a final remark, we stress the fact that the proposed approach supports and complements the remarkable step taken by the ATLAS collaboration in publishing full experimental likelihoods [12, 14, 15]. On one hand, the DNNLikelihood framework can be seen as a way to export a HistFactory likelihood to a software environment that doesn't meet the needed dependencies, or as an alternative option for experimental collaborations that don't use HistFactory. On the other, the use of a DNN model doesn't imply specific choices on the likelihood function (e.g. it covers equally well unbinned likelihoods and/or likelihoods built from analytical functions and not represented as histograms).

Last but not least, using a DNNLikelihood could be a viable solution to distribute the outcome of phenomenological studies, such as the ones presented in refs. [19–28, 67, 68].

On a long term, a large scale adoption of likelihood publishing towards the DNNLikelihood framework would motivate the possibility of publishing DNN models on HEPData

(which, more generically, would be beneficial for reproducibility issues related to physics analysis using DNNs for selection, etc.), for instance supporting the ONNX format. In this respect, and considering that there could be interest in likelihood publishing in other domains (e.g. for the Λ_{CDM} likelihood in cosmology), it might be worth considering a less hierarchical submission procedure for HEPData (e.g. allowing individuals outside a structured organization/experimental collaboration to submit a likelihood function) or the opportunity to create a separate (and not HEP specific) likelihood function repository, e.g. based on Zenodo.

Acknowledgements

We are indebted to the authors of Ref. [9] for providing us with details and data concerning the LHC-like experiment considered in this paper. This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (grant agreement n° 772369) and from the Italian Ministry of Research (MIUR) under grant PRIN 20172LNEEZ. The work of R. Torre has been partially supported by the INFN Grant MALHEPHYCA.

A On the multivariate normal distribution

For the univariate normal distribution (centered in zero and with unit variance)

$$\mathcal{N}_1(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}, \quad (19)$$

the confidence intervals at $n\sigma$ are defined by the quantiles of the χ_k^2 distribution with $k = 1$ degrees of freedom by

$$n^2 = 2Q^{-1}\left(\frac{1}{2}, 0, 1 - \alpha\right), \quad \alpha = Q\left(\frac{1}{2}, 0, \frac{n^2}{2}\right) \quad (20)$$

where $1 - \alpha$ is the area of the distribution within $\pm n\sigma$ from the mean. Generalization to a multivariate normal distribution in l dimensions with zero vector mean and identity matrix variance $\Sigma = I_l$

$$\mathcal{N}_l(\mathbf{x}) = \frac{1}{(2\pi)^{l/2}} e^{-\frac{1}{2}x_i \Sigma^{ij} x_j}, \quad (21)$$

is obtained by considering the χ_k^2 distribution with $k = l$ degrees of freedom. Equation (20) then becomes

$$n^2 = 2Q^{-1}\left(\frac{l}{2}, 0, \alpha\right), \quad \alpha = Q\left(\frac{l}{2}, 0, \frac{n^2}{2}\right). \quad (22)$$

The typical size of N_l outside of the confidence interval at $n\sigma$ is given by

$$\mathcal{N}_l(\mathbf{x} \sim \|\mu - n\sigma\|) \approx \frac{1}{(2\pi)^{l/2}} e^{-\frac{1}{2}2Q^{-1}\left(\frac{l}{2}, 0, 1 - \alpha\right)}, \quad (23)$$

where we have denoted by $\mathbf{x} \sim \|\mu - n\sigma\|$ a point that is approximately $n\sigma$ away from the l dimensional mean. The target set of the distribution within $n\sigma$ from the mean is therefore given by:

$$\mathcal{N}_l(\mathbf{x} \lesssim \|\mu - n\sigma\|) \in (2\pi)^{-l/2} \{\min(\mathcal{N}_l(\mathbf{x}))|_{\mathbf{x} \lesssim \|\mu - n\sigma\|}, 1\} = \{e^{-\frac{1}{2}2Q^{-1}(\frac{l}{2}, 0, 1-\alpha)}, 1\}, \quad (24)$$

where $\min(\mathcal{N}_l(\mathbf{x}))|_{\mathbf{x} \lesssim \|\mu - n\sigma\|}$ is the minimum of \mathcal{N}_l for \mathbf{x} within $n\sigma$ from the mean divided by the normalisation factor $(2\pi)^{-l/2}$. In Table 6 we show the value of $-\log_{10}(\min(\mathcal{N}_l(\mathbf{x}))|_{\mathbf{x} \lesssim \|\mu - n\sigma\|})$ for $l = 1, 10, 100, 1000$ dimensions and for confidence intervals up to 6σ .

$n \setminus l$	1	10	100	1000
1	0	3	23	222
2	1	4	27	234
3	2	6	31	245
4	3	8	36	256
5	5	10	40	268
6	8	13	45	279

Table 6: Value of $-\log_{10}(\min(\mathcal{N}_l(\mathbf{x}))|_{\mathbf{x} \lesssim \|\mu - n\sigma\|})$ for $l = 1, 10, 100, 1000$ dimensions and for confidence intervals up to 6σ . See the text for details.

B Pseudo-experiments and frequentist coverage

In order to check the validity of the asymptotic approximation given by Wilks’ theorem, and to obtain more robust upper bounds from the t_μ test-statistics for the LHC-like search discussed in Section 3, we need to generate several pseudo-experiments for each hypothesised “true” value of μ , compute the log-likelihood and the test statistics t_μ for each pseudo-experiment, and study the pdf of t_μ , denoted by $f(t_\mu|\mu)$. The cumulative distribution of $f(t_\mu|\mu)$ determines the coverage properties of the test-statistics t_μ through the relation

$$1 - p_\mu = \int_0^{t_{\mu, \text{obs}}} f(t_\mu|\mu) dt_\mu, \quad (25)$$

where $t_{\mu, \text{obs}}$ is the t_μ of our actual experiment at the given value of μ . By solving this equation for μ at specified value of $p_\mu = \alpha$ we get the value of μ excluded at a CL of $1 - \alpha$.

We generate pseudo-experiments following two different procedures for the treatment of nuisance parameters and compare the results.

- **Profile construction**

First we consider a frequentist treatment of nuisance parameters, referred to as profile construction [41, 69]: we determine the values $\hat{\theta}_\mu$ of the nuisance parameters at the maximum of the LF of our actual experiment for different values of μ , and keep them fixed

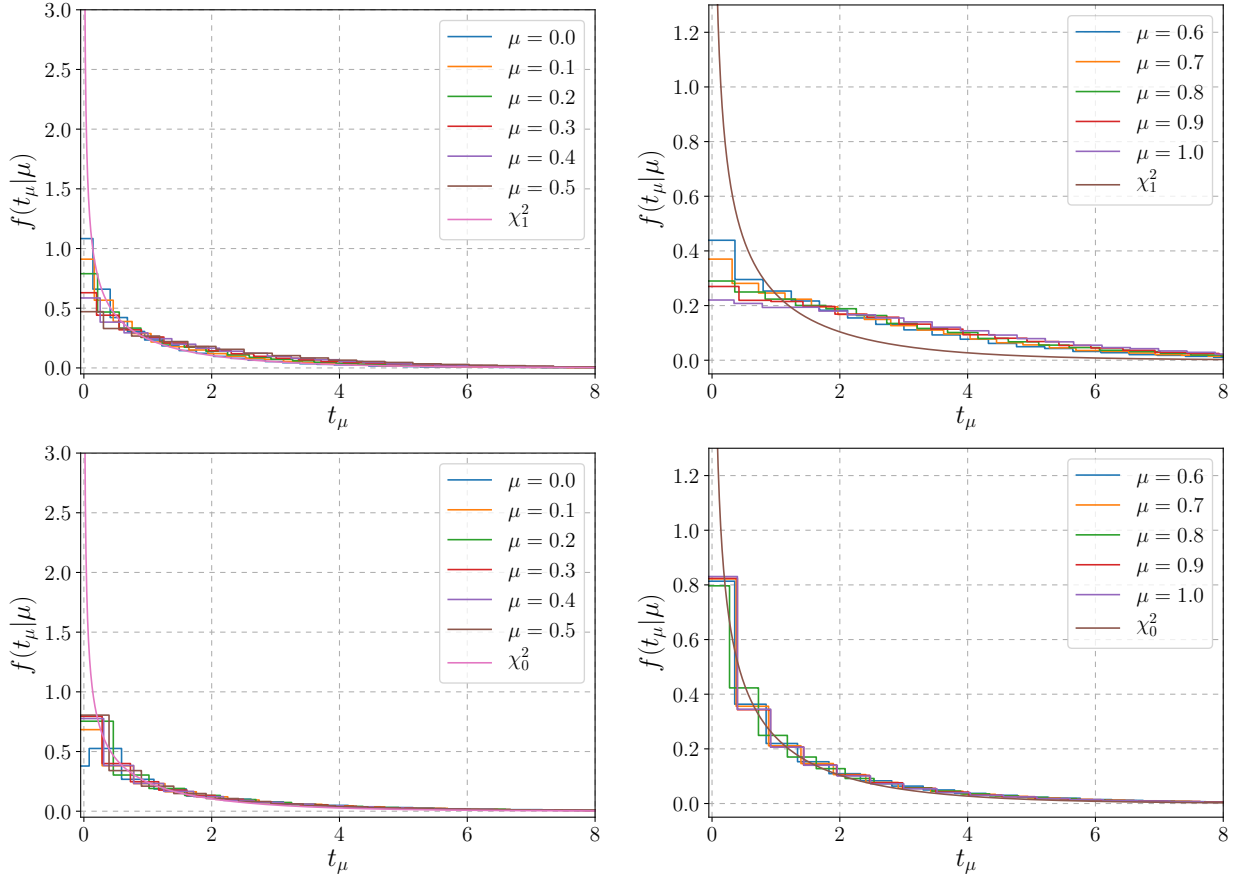


Figure 17: Distributions $f(t_\mu|\mu)$ for the LHC-like new physics search discussed in Section 3, obtained from pseudo-experiments employing the profile construction (upper panel) and the hybrid frequentist-Bayesian treatment of nuisance parameters (lower panel). A χ^2_1 distribution is shown for comparison.

to generate several pseudo-experiments for each μ according to eq. (12). This approach encodes statistical fluctuations arising in repeated experiments, but does not take into account systematic uncertainties. Although this approach violates the “anticipation criterion” [70], it is expected to work well, and have good coverage properties when the uncertainty is statistically dominated.¹³ Indeed the coverage of this approach grows as the profiled value $\hat{\theta}_\mu$ approaches the “true” value of θ . This happens when systematic uncertainties become less and less relevant compared to statistical fluctuations.

Following this procedure we generated $5 \cdot 10^4$ pseudo-experiments for each value of μ in the range $[0, 1]$ with steps of 0.1.

¹³This was the approach employed, for instance, in the LHC Higgs boson search combination in 2011 [3].

μ	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
$t_{\mu,\text{obs}}$	0.0	0.09	0.29	0.63	1.13	1.77	2.57	3.53	4.64	5.92	7.37
coverage (profile)	0.0	0.15	0.27	0.37	0.49	0.59	0.68	0.76	0.83	0.89	0.93
coverage (hybrid)	0.0	0.15	0.30	0.47	0.61	0.73	0.83	0.90	0.94	0.97	0.99
coverage (χ_1^2)	0.0	0.23	0.41	0.57	0.71	0.82	0.89	0.94	0.97	0.99	0.99

Table 7: Values of μ with corresponding $t_{\mu,\text{obs}}$ and CL coverage for the LHC-like new physics search discussed in Section 3, obtained from pseudo-experiments employing the profile construction and the hybrid frequentist-Bayesian treatment of nuisance parameters. The last row shows the asymptotic result obtained with a χ_1^2 distribution.

CL	68.27%	95.45%	99.73%
t_μ (profile)	2.71	> 7.37	> 7.37
μ (profile)	0.61	> 1.00	> 1.00
t_μ (hybrid)	1.52	5.26	> 7.37
μ (hybrid)	0.46	0.85	> 1.00
t_μ (χ_1^2)	1.00	4.00	9.00
μ (χ_1^2)	0.37	0.74	> 1.00

Table 8: Upper bounds on t_μ and corresponding μ at different CL for the LHC-like new physics search discussed in Section 3, obtained from pseudo-experiment employing a profiled frequentist and hybrid frequentist-Bayesian treatment of nuisance parameters. The last two rows show the asymptotic result obtained with a χ_1^2 distribution.

- **Hybrid frequentist-Bayesian (marginal model)**

In order to understand the impact of systematic uncertainties we also consider the hybrid frequentist-Bayesian treatment of nuisance parameters (referred to as marginal model in ref. [41]): pseudo-experiments are generated including both variations of the nuisance parameters according to their distribution, and statistical uncertainty through eq. (12). This allows to include the effect of systematic uncertainties in the generation of pseudo-experiments.

In this case we generated, for the same values of μ considered above, 10^4 expected counts (for all 90 bins) corresponding to variations of the nuisance parameters over their multivariate distribution, and subsequently used each of these expected counts to generate 10 pseudo-experiments according to eq. (12). This delivers a total of 10^5 pseudo-experiments for each value of μ , encoding both statistical and systematic uncertainties.

Figure 17 shows the distributions $f(t_\mu|\mu)$ for each value of μ , while the probability values covered for each value of μ and the corresponding $t_{\mu,\text{obs}}$ are given in Table 7. Using the distributions $f(t_\mu|\mu)$ we performed the integral in eq. (25) for $1 - p_\mu = 0.6827, 0.9545, 0.9973$, getting the upper bounds on μ reported in Table 8. All results are shown for both the profile

construction and the hybrid frequentist-Bayesian approach.

As expected, due to the small statistics in several bins, the asymptotic result is inaccurate, and in particular tends to undercover the true value. This delivers a too “aggressive” upper bound for μ . The hybrid approach gives relatively more conservative upper bounds, very similar (expectedly, given the marginal model used for the nuisance parameters) to the Bayesian result reported in Table 1. Finally, the profile construction gives the most conservative bound, that is substantially more conservative than the asymptotic one. We do not dare here to discuss which upper bound is to be quoted, since we are only interested in assessing the validity of the asymptotic approximation.

References

- [1] A. Stuart, J. K. Ord and S. Arnold, *Kendall’s advanced theory of statistics. Vol.2A: Classical inference and the linear model (Sixth Edition)*, John Wiley & Sons, Ltd, 2009 [CDS].
- [2] A. O’Hagan and J. Forster, *Kendall’s advanced theory of statistics. Vol.2B: Bayesian inference (Second Edition)*, John Wiley & Sons, Ltd, 2004 [CDS].
- [3] ATLAS, CMS and LHC HIGGS COMBINATION GROUP Collaborations, *Procedure for the LHC Higgs boson search combination in Summer 2011*, Tech. Rep. CMS-NOTE-2011-005, ATL-PHYS-PUB-2011-11, CERN, 2011 [INSPIRE].
- [4] FCC Collaboration, A. Abada et al., *FCC Physics Opportunities*, *Eur. Phys. J.* **C79** (2019) 474 [INSPIRE].
- [5] T. Behnke et al., *The International Linear Collider Technical Design Report - Volume 1: Executive Summary*, 1306.6327 [INSPIRE].
- [6] M. Aicheler, P. Burrows, M. Draper, T. Garvey, P. Lebrun, K. Peach et al., *A Multi-TeV Linear Collider Based on CLIC Technology - CLIC Conceptual Design Report*, CERN Yellow Reports: Monographs CERN, Geneva, 2012 [CDS].
- [7] N. Berger et al., *Simplified Template Cross Sections - Stage 1.1*, 1906.02754 [INSPIRE].
- [8] S. Fichtel, *Taming systematic uncertainties at the LHC with the central limit theorem*, *Nucl. Phys.* **B911** (2016) 623 [1603.03061] [INSPIRE].
- [9] A. Buckley, M. Citron, S. Fichtel, S. Kraml, W. Waltenberger and N. Wardle, *The Simplified Likelihood Framework*, *JHEP* **04** (2019) 064 [1809.05548] [INSPIRE].
- [10] CMS Collaboration, *Simplified likelihood for the re-interpretation of public CMS results*, Tech. Rep. CMS-NOTE-2017-001, 2017 [CDS].

- [11] K. Cranmer, S. Kreiss, D. Lopez-Val and T. Plehn, *Decoupling Theoretical Uncertainties from Measurements of the Higgs Boson*, *Phys. Rev.* **D91** (2015) 054032 [[1401.0080](#)] [[INSPIRE](#)].
- [12] ATLAS Collaboration, G. Aad et al., *Search for bottom-squark pair production with the ATLAS detector in final states containing Higgs bosons, b-jets and missing transverse momentum*, [1908.03122](#) [[INSPIRE](#)].
- [13] G. Watt et al., *HepData* [[WEBPAGE](#)].
- [14] ROOT Collaboration, K. Cranmer et al., *HistFactory: A tool for creating statistical models for use with RooFit and RooStats*, Tech. Rep. CERN-OPEN-2012-016, 2012 [[CDS](#)].
- [15] ATLAS Collaboration, *Reproducing searches for new physics with the ATLAS experiment through publication of full statistical likelihoods*, Tech. Rep. ATL-PHYS-PUB-2019-029, 2019 [[CDS](#)].
- [16] F. James, Y. Perrin and L. Lyons, eds., *1st Workshop on Confidence Limits, CERN, Geneva, Switzerland, 17-18 Jan 2000: Proceedings*, *CERN*, 2000 [[INSPIRE](#)].
- [17] CMS Collaboration, A. M. Sirunyan et al., *Measurements of properties of the Higgs boson decaying into the four-lepton final state in pp collisions at $\sqrt{s} = 13$ TeV*, *JHEP* **11** (2017) 047 [[1706.09936](#)] [[INSPIRE](#)].
- [18] J. Bai et al., *ONNX: Open Neural Network Exchange* [[GITHUB](#)].
- [19] M. Ciuchini, G. D'Agostini, E. Franco, V. Lubicz, G. Martinelli, F. Parodi et al., *2000 CKM triangle analysis: A Critical review with updated experimental inputs and theoretical parameters*, *JHEP* **07** (2001) 013 [[hep-ph/0012308](#)] [[INSPIRE](#)].
- [20] A. Hocker, H. Lacker, S. Laplace and F. Le Diberder, *A New approach to a global fit of the CKM matrix*, *Eur. Phys. J.* **C21** (2001) 225 [[hep-ph/0104062](#)] [[INSPIRE](#)].
- [21] J. Charles, A. Hocker, H. Lacker, S. Laplace, F. R. Le Diberder, J. Malcles et al., *CP violation and the CKM matrix: Assessing the impact of the asymmetric B factories*, *Eur. Phys. J.* **C41** (2005) 1 [[hep-ph/0406184](#)] [[INSPIRE](#)].
- [22] UTFIT Collaboration, M. Bona et al., *The 2004 UTfit collaboration report on the status of the unitarity triangle in the standard model*, *JHEP* **07** (2005) 028 [[hep-ph/0501199](#)] [[INSPIRE](#)].
- [23] UTFIT Collaboration, M. Bona et al., *Model-independent constraints on $\Delta F = 2$ operators and the scale of new physics*, *JHEP* **03** (2008) 049 [[0707.0636](#)] [[INSPIRE](#)].

- [24] M. Ciuchini, E. Franco, S. Mishima and L. Silvestrini, *Electroweak Precision Observables, New Physics and the Nature of a 126 GeV Higgs Boson*, *JHEP* **08** (2013) 106 [[1306.4644](#)] [[INSPIRE](#)].
- [25] GFITTER GROUP Collaboration, M. Baak et al., *The global electroweak fit at NNLO and prospects for the LHC and ILC*, *Eur. Phys. J.* **C74** (2014) 3046 [[1407.3792](#)] [[INSPIRE](#)].
- [26] J. de Blas, M. Ciuchini, E. Franco, S. Mishima, M. Pierini, L. Reina et al., *Electroweak precision observables and Higgs-boson signal strengths in the Standard Model and beyond: present and future*, *JHEP* **12** (2016) 135 [[1608.01509](#)] [[INSPIRE](#)].
- [27] A. Falkowski, M. González-Alonso and K. Mimouni, *Compilation of low-energy constraints on 4-fermion operators in the SMEFT*, *JHEP* **08** (2017) 123 [[1706.03783](#)] [[INSPIRE](#)].
- [28] J. Ellis, C. W. Murphy, V. Sanz and T. You, *Updated Global SMEFT Fit to Higgs, Diboson and Electroweak Data*, *JHEP* **06** (2018) 146 [[1803.03252](#)] [[INSPIRE](#)].
- [29] M. Clark, *MCMC Algorithms* [[WEBPAGE](#)].
- [30] GEANT4 Collaboration, S. Agostinelli et al., *GEANT4: A Simulation toolkit*, *Nucl.Instrum.Meth.A* **506** (2003) 250 [[INSPIRE](#)].
- [31] K. Kandasamy, J. Schneider and B. Póczos, *Query Efficient Posterior Estimation in Scientific Experiments via Bayesian Active Learning*, *Artificial Intelligence* **243** (2017) 45 [[1702.01145](#)] [[SEMANTIC SCHOLAR](#)].
- [32] S. Caron, T. Heskes, S. Otten and B. Stienen, *Constraining the Parameters of High-Dimensional Models with Active Learning*, [1905.08628](#) [[INSPIRE](#)].
- [33] A. Coccaro, M. Pierini, L. Silvestrini and R. Torre, *to appear*.
- [34] F. Feroz, K. Cranmer, M. Hobson, R. Ruiz de Austri and R. Trotta, *Challenges of Profile Likelihood Evaluation in Multi-Dimensional SUSY Scans*, *JHEP* **06** (2011) 042 [[1101.3296](#)] [[INSPIRE](#)].
- [35] S. Kullback and R. A. Leibler, *On Information and Sufficiency*, *Ann. Math. Statist.* **22** (1951) 79 [[SEMANTIC SCHOLAR](#)].
- [36] A. Kolmogorov, *Sulla Determinazione Empirica di una Legge di Distribuzione*, *Giornale dell'Istituto Italiano degli Attuari* **4** (1933) 83 [[GOOGLE SCHOLAR](#)].
- [37] N. Smirnov, *Table for Estimating the Goodness of Fit of Empirical Distributions*, *Ann. Math. Statist.* **19** (1948) 279 [[SEMANTIC SCHOLAR](#)].

- [38] B. Krawczyk, *Learning from imbalanced data: open challenges and future directions*, *Progress in Artificial Intelligence* **5** (2016) 221 [[SEMANTIC SCHOLAR](#)].
- [39] P. Branco, L. Torgo and R. P. Ribeiro, *SMOBN: a Pre-processing Approach for Imbalanced Regression*, in *First International Workshop on Learning with Imbalanced Domains: Theory and Applications, LIDTA@PKDD/ECML 2017, 22 September 2017, Skopje, Macedonia*, vol. 74 of *Proceedings of Machine Learning Research*, pp. 36–50, PMLR, 2017 [[SEMANTIC SCHOLAR](#)].
- [40] M. Ren, W. Zeng, B. Yang and R. Urtasun, *Learning to Reweight Examples for Robust Deep Learning*, in *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, pp. 4331–4340, 2018 [[1803.09050](#)] [[SEMANTIC SCHOLAR](#)].
- [41] PARTICLE DATA GROUP Collaboration, M. Tanabashi et al., *Review of Particle Physics*, *Phys. Rev.* **D98** (2018) 030001 [[INSPIRE](#)].
- [42] D. Foreman-Mackey, D. W. Hogg, D. Lang and J. Goodman, *emcee: The MCMC Hammer*, *Publ. Astron. Soc. Pac.* **125** (2013) 306 [[1202.3665](#)] [[INSPIRE](#)].
- [43] J. Goodman and J. Weare, *Ensemble samplers with affine invariance*, *Communications in Applied Mathematics and Computational Science, Vol. 5, No. 1, p. 65-80, 2010* **5** (2010) 65 [[SEMANTIC SCHOLAR](#)].
- [44] A. Gelman and D. B. Rubin, *Inference from Iterative Simulation Using Multiple Sequences*, *Statist. Sci.* **7** (1992) 457 [[INSPIRE](#)].
- [45] S. P. Brooks and A. Gelman, *General Methods for Monitoring Convergence of Iterative Simulations*, *Journal of Computational and Graphical Statistics* **7** (1998) 434 [[SEMANTIC SCHOLAR](#)].
- [46] D. Huijser, J. Goodman and B. J. Brewer, *Properties of the Affine Invariant Ensemble Sampler in high dimensions*, [1509.02230](#) [[SEMANTIC SCHOLAR](#)].
- [47] D. Foreman-Mackey, D. W. Hogg, D. Lang and J. Goodman, *emcee: The MCMC Hammer* [[READTHEDOCS](#)].
- [48] D. Foreman-Mackey, D. W. Hogg, D. Lang and J. Goodman, *emcee: The MCMC Hammer* [[GITHUB](#)].
- [49] E. B. Ford, *Convergence Diagnostics For Markov chain Monte Carlo*, 2016 [[SLIDES](#)].
- [50] W. A. Link and M. J. Eaton, *On thinning of chains in MCMC*, *Methods in Ecology and Evolution* **3** (2012) 112 [[SEMANTIC SCHOLAR](#)].
- [51] A. B. Owen, *Statistically efficient thinning of a Markov chain sampler*, [1510.07727](#) [[SEMANTIC SCHOLAR](#)].

- [52] A. Cocco, M. Pierini, L. Silvestrini and R. Torre, *to appear*.
- [53] S. S. Wilks, *The Large-Sample Distribution of the Likelihood Ratio for Testing Composite Hypotheses*, *Annals Math. Statist.* **9** (1938) 60 [INSPIRE].
- [54] G. Cowan, K. Cranmer, E. Gross and O. Vitells, *Asymptotic formulae for likelihood-based tests of new physics*, *Eur. Phys. J.* **.17** (2011) 1554 [1007.1727] [Erratum: *Eur. Phys. J.* **C73** (2013) 2501] [INSPIRE].
- [55] F. Chollet et al., *Keras: Deep Learning for humans* [GITHUB].
- [56] M. Abadi et al., *TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems* [TensorFlow v1].
- [57] C. Zhang, S. Bengio, M. Hardt, B. Recht and O. Vinyals, *Understanding deep learning requires rethinking generalization*, in *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017 [1611.03530] [SEMANTIC SCHOLAR].
- [58] M. Belkin, D. Hsu and P. Mitra, *Overfitting or perfect fitting? Risk bounds for classification and regression rules that interpolate*, in *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada*, pp. 2306–2317, 2018 [1806.05161] [SEMANTIC SCHOLAR].
- [59] X. Glorot, A. Bordes and Y. Bengio, *Deep Sparse Rectifier Neural Networks*, in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, G. Gordon, D. Dunson and M. Dudík, eds., vol. 15 of *Proceedings of Machine Learning Research*, pp. 315–323, PMLR, 2011 [SEMANTIC SCHOLAR].
- [60] D.-A. Clevert, T. Unterthiner and S. Hochreiter, *Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs)*, in *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, Nov, 2016 [1511.07289] [SEMANTIC SCHOLAR].
- [61] G. Klambauer, T. Unterthiner, A. Mayr and S. Hochreiter, *Self-Normalizing Neural Networks*, in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pp. 971–980, 2017 [1706.02515] [SEMANTIC SCHOLAR].
- [62] Y. LeCun, L. Bottou, G. B. Orr and K.-R. Müller, *Efficient BackProp*, in *Neural Networks: Tricks of the Trade. Lecture Notes in Computer Science*, M. G., O. G.B. and M. KR., eds., vol. 7700, Springer-Verlag, 2012 DOI [SEMANTIC SCHOLAR].

- [63] S. L. Smith, P.-J. Kindermans, C. Ying and Q. V. Le, *Don't Decay the Learning Rate, Increase the Batch Size*, in *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*, 2018 [[1711.00489](#)] [[SEMANTIC SCHOLAR](#)].
- [64] Y. Yao, L. Rosasco and A. Caponnetto, *On Early Stopping in Gradient Descent Learning*, *Constructive Approximation* **26** (2007) 289 [[SEMANTIC SCHOLAR](#)].
- [65] G. Raskutti, M. J. Wainwright and B. Yu, *Early stopping and non-parametric regression: an optimal data-dependent stopping rule*, *The Journal of Machine Learning Research* **15** (2014) 335 [[1306.3574](#)] [[SEMANTIC SCHOLAR](#)].
- [66] L. Luo, Y. Xiong, Y. Liu and X. Sun, *Adaptive Gradient Methods with Dynamic Bound of Learning Rate*, in *Proceedings of the 7th International Conference on Learning Representations*, (New Orleans, Louisiana), May, 2019 [[1902.09843](#)] [[SEMANTIC SCHOLAR](#)].
- [67] M. Ciuchini, A. M. Coutinho, M. Fedele, E. Franco, A. Paul, L. Silvestrini et al., *New Physics in $b \rightarrow s\ell^+\ell^-$ confronts new data on Lepton Universality*, *Eur. Phys. J.* **C79** (2019) 719 [[1903.09632](#)] [[INSPIRE](#)].
- [68] I. Brivio, S. Bruggisser, F. Maltoni, R. Moutafis, T. Plehn, E. Vryonidou et al., *O new physics, where art thou? A global search in the top sector*, [1910.03606](#) [[INSPIRE](#)].
- [69] K. Cranmer, *Statistical challenges for searches for new physics at the LHC*, in *Statistical Problems in Particle Physics, Astrophysics and Cosmology (PHYSTAT 05): Proceedings, Oxford, UK, September 12-15, 2005*, pp. 112–123, 2005 [[physics/0511028](#)], DOI [[INSPIRE](#)].
- [70] L. Demortier, *Constructing Ensembles of Pseudo-Experiments*, *eConf* **C030908** (2003) WEMT003 [[physics/0312100](#)] [[INSPIRE](#)].