

Aalto University  
School of Science  
Master's Programme in Life Science Technologies

Bent Ivan Oliver Harnist

# Probabilistic Precipitation Nowcasting using Bayesian Convolutional Neural Networks

Master's Thesis  
Espoo, July 20th, 2022

**DRAFT! — July 21, 2022 — DRAFT!**

Supervisor: Professor Arno Solin  
Advisor: Terhi Mäkinen D.Sc.  
Seppo Pulkkinen D.Sc.

<b>Author:</b>	Bent Ivan Oliver Harnist	
<b>Title:</b>		
Probabilistic Precipitation Nowcasting using Bayesian Convolutional Neural Networks		
<b>Date:</b>	July 20th, 2022	<b>Pages:</b> 87
<b>Major:</b>	Complex Systems	<b>Code:</b> SCI3060
<b>Supervisor:</b>	Professor Arno Solin	
<b>Advisor:</b>	Terhi Mäkinen D.Sc. Seppo Pulkkinen D.Sc.	
<p>Precipitation nowcasting refers to forecasting precipitation at timescales of dozens of minutes to a few hours, and is usually approached using models based on real-time radar images of precipitation. Now, reliably nowcasting precipitation probabilities locally is a societally important challenge, due to the danger and economic losses potentially incurred by unexpected precipitation, such as the damages caused by large hail or flash floods due to heavy localized rainfall.</p> <p>In this thesis, a Convolutional Neural Network (CNN) used for nowcasting precipitation point estimates is adapted into a Bayesian Neural Network (BNN) for producing ensemble predictions and estimate precipitation probabilities. A BNN turns model parameter values into probability distributions that are estimated using bayesian inference. In this work, parameters are modeled using gaussian distributions, whose parameters are learned through Stochastic Variational Inference. The predictions are then simply formed by Monte-Carlo sampling each parameter from its distribution, inducing a predictive distribution.</p> <p>The developed model, titled BCNN, has its predictive skill evaluated against several baseline models using a multitude of criteria covering various aspects of the nowcasts. On some of those aspects, the model was found to perform at a skill close to or equal to baselines. However, it was also found to be challenging to produce a well-calibrated and reliable probabilistic nowcast for BCNN, owing to likely properties of the model and problems in the training procedure. Nevertheless, BCNN provided encouraging preliminary results, paving the way to more solid attempts at probabilistic precipitation nowcasting using Bayesian Neural Networks. Particularly, a more thoughtful modeling and integration of the data-related uncertainties may make the nowcasts more reliable by making the predictive distribution less sensitive to initial conditions.</p>		
<b>Keywords:</b>	Precipitation nowcasting,	
<b>Language:</b>	English	

Aalto-yliopisto

Perustieteiden korkeakoulu

Tieto-, tietoliikenne- ja informaatiotekniikan maisteriohjelma

DIPLOMITYÖN

TIIVISTELMÄ

<b>Tekijä:</b>	Bent Ivan Oliver Harnist		
<b>Työn nimi:</b>	Probabilistinen Sateen Nowcasting käyttäen Bayesilaisia Konvolutiivisia Neuroverkkoja		
<b>Päiväys:</b>	20. heinäkuuta 2022	<b>Sivumäärä:</b>	87
<b>Pääaine:</b>	Complex Systems	<b>Koodi:</b>	SCI3060
<b>Valvoja:</b>	Professori Arno Solin		
<b>Ohjaaja:</b>	Terhi Mäkinen D.Sc. Seppo Pulkkinen D.Sc.		
Tiivistelmä			
<b>Asiasanat:</b>	Sateen ennustaminen,		
<b>Kieli:</b>	Englanti		

# Acknowledgements

**!FIXME My acknowledgements FIXME!**

Espoo, July 20th, 2022

Bent Ivan Oliver Harnist

# Abbreviations and Acronyms

DL	Deep learning
NWP	Numerical Weather Prediction
BNN	Bayesian Neural Network
BCNN	Bayesian Convolutional Neural Network
STEPS	Short-Term Ensemble Prediction System
LINDA	Lagrangian Integro-Difference equation model with Autoregression
Z	Reflectivity factor
dBZ	Decibel relative to Z
ELBO	Evidence Lower Bound
CSI	Critical Success Index
ETS	Equitable Threat Score
FAR	False Alarm Rate
POD	Probability of Detection
MAE	Mean Absolute Error
ME	Mean Error
FSS	Fractions Skill Score
RAPSD	Radially-averaged Power Spectral Density
CRPS	Continuous Ranked Probability Score
ROC	Receiver Operating Characteristic
WMO	World Meteorological Organization

# Contents

<b>Abbreviations and Acronyms</b>	<b>5</b>
<b>1 Introduction</b>	<b>8</b>
1.1 Problem statement . . . . .	9
1.2 Structure of the Thesis . . . . .	10
<b>2 Background</b>	<b>11</b>
2.1 Precipitation Nowcasting . . . . .	11
2.1.1 Classical deterministic methods . . . . .	14
2.1.2 Classical probabilistic methods . . . . .	16
2.2 Deep Learning for Nowcasting . . . . .	17
2.2.1 Artificial Neural Networks . . . . .	17
2.2.2 Deep Learning approaches to Nowcasting . . . . .	19
2.3 Bayesian Deep Learning . . . . .	21
2.3.1 Learning Probability Distributions . . . . .	21
2.3.2 Variational inference . . . . .	23
2.3.3 Predictive uncertainty estimation and decomposition .	26
<b>3 Materials and Methods</b>	<b>29</b>
3.1 Experimental Setup . . . . .	29
3.1.1 Software and resources . . . . .	31
3.2 Dataset and data selection . . . . .	32
3.3 Model . . . . .	33
3.3.1 The baseline: A U-Net based CNN . . . . .	33
3.3.2 BCNN: a Bayesian extension to RainNet . . . . .	35
3.3.3 Additional Hyper-parameters for NN . . . . .	38
3.4 Verification Baseline Models . . . . .	39
3.4.1 Baseline Deterministic Models . . . . .	39
3.4.2 Baseline Probabilistic Models . . . . .	39
3.5 Verification Visualizations and Metrics . . . . .	40
3.5.1 Deterministic Prediction skill evaluation metrics . . . . .	40

3.5.2	Visual Evaluation of nowcast predictive uncertainty and Skill . . . . .	43
3.5.3	Probabilistic Prediction skill evaluation metrics . . . . .	43
<b>4</b>	<b>Results</b>	<b>46</b>
4.1	BCNN Hyperparameter Optimization . . . . .	46
4.2	Nowcasting Examples . . . . .	49
4.2.1	Case study 1 : Rapidly Moving Mixed Stratiform and Convective Rain Event . . . . .	49
4.2.2	Case study 2 : Mixed Stratiform and Convective Rain Event with Spinning Motion . . . . .	51
4.2.3	Additional Case Study 1 results . . . . .	53
4.3	Deterministic prediction skill (Metrics) . . . . .	54
4.4	Probabilistic prediction skill (Metrics) . . . . .	57
<b>5</b>	<b>Discussion</b>	<b>62</b>
5.1	Goodness of results . . . . .	62
5.2	Validity of Assumptions and Experiments . . . . .	63
5.3	Directions for Further Work . . . . .	66
<b>6</b>	<b>Conclusions</b>	<b>68</b>
<b>7</b>	<b>references</b>	<b>70</b>
<b>A</b>	<b>Additional Figures</b>	<b>78</b>
A.1	Model Selection . . . . .	79
A.2	Case Studies . . . . .	81
A.3	Deterministic Metrics . . . . .	86
A.4	Probabilistic Metrics . . . . .	87

# Chapter 1

## Introduction

Nowcasting is defined as weather forecasting at a local scale up to six hours, according the World Meteorological Organization (WMO) definition [54]. The ability to provide an accurate precipitation nowcast has grown during this century into a meteorologically and societally significant challenge. This significance emanates from the fact that early warnings of severe precipitation enable authorities and other actors to make early decisions enabling for example disaster damage control, traffic safety, and in the case of heavy rainfall flash flood prevention, as well as the mitigation of economic loss incurred by them. High-density urbanization has exacerbated these issues, making it evermore vital to discover reliable and skillful methods for precipitation nowcasting.

Numerical weather prediction (NWP) solves the partial differential equations governing the physical processes of weather and climate to produce forecasts. NWP has seen great improvements over decades, up to the point where it can be used to produce accurate forecasts for up to a week in some cases [9]. However, such method is not applicable to predicting at timescales as short as 0 to 6 hours. This is mainly due to imperfect initialization making simulations unable to reach numerical stability in such short time timescales. Other problems with NWP include its computational expensiveness and generally having insufficient spatiotemporal to predict convective events and heavy localized rainfall [56].

To compensate for the shortcomings of NWP in the realm of precipitation nowcasting, many different dedicated models and algorithms have been developed [42]. Contrary to NWP, these models do not combine data from multiple sources in making predictions, but are simpler in a way as they use only a single input field. In practice, many of those systems are based on the estimation of the precipitation advection field from radar echo image sequences and the further extrapolation of these sequences along the advec-

tion field [42, 47]. Other methods have been developed that track and try to nowcast the evolution of individual rain cells instead of the whole grid [18, 42].

Dedicated nowcasting methods have had great success in forecasting the immediate future, but their performance typically degrades quickly, becoming unreliable often in the range of an hour. This is related to the fact that basic extrapolation-based methods have limitations such as failing to capture nonlinear patterns like growth and decay of precipitation, as well as convective cell initiation and their life-cycle.

A lot of work has been done in nowcasting in trying to incorporate modeling of higher-order and multiscale phenomena into advection-extrapolation based models. Recently, Deep-learning (DL) based precipitation nowcasting approaches have started showing promising results delving into the issues of traditional methods thanks to the high volume of training data available, the increase in computational power, the expressivity of those models and their relative cheapness of inference [8, 59, 60].

## 1.1 Problem statement

While the above introduction has cast the problem of nowcasting as finding a single optimal point estimate for future precipitation, it is useful to think about it in a probabilistic way. Because in nowcasting the most interesting phenomena are extreme events, which are those requiring preparation and early warnings, it makes intuitive sense that accurate and reliable probabilistic nowcasts are primordial for operational decision-making in meteorological crises. Also, because smaller spatial scales are less predictable, Adding a degree of uncertainty improves the model usefulness at those scales [24].

Indeed, the interest in probabilistic nowcasts has grown lately and many probabilistic models have been introduced in the last decades [5, 21, 55, 58], notably the Short-Term Ensemble Prediction System (STEPS) [13] and more recently the Lagrangian Integro-Difference equation model with Autoregression (LINDA) [43]. These methods predict ensembles, that are used to estimate underlying distributions of data and precipitation probabilities. They perform reasonably well, but suffer from the same limitations as other extrapolation based methods, namely a limited ability to predict nonlinear patterns of growth and decay.

So far, only little work has been done on using DL to produce probabilistic precipitation nowcasts, with the biggest breakthrough perhaps being Ravuri et al. [46] using adversarially trained deep generative models to produce ensemble nowcasts. There exists many possible ways of making probabilistic

nowcasts using deep learning, with most of them focusing on directly modeling probability distributions of data. Another approach, that will be focused on from now on is instead to model the uncertainty of model parameters rather than that of the data, as would be done with STEPS or LINDA. This is done by using Bayesian Neural Networks and has an advantage of providing implicit regularization of model parameters thus reducing overfitting, while outputting an ensemble of predictions whose variability in part reflects network parameter uncertainty given the training data. Bayesian neural networks have been proposed for use in other risk-averse application such as biomedical image segmentation [32] and autonomous vehicles [37].

In this work, the problem of making skillful and reliable probabilistic precipitation nowcasts will approached by building an uncertainty-aware Neural Network. A baseline Convolutional Neural Network (CNN) will be turned into a Bayesian Convolutional Network (BCNN) with optimization performed with stochastic variational inference (SVI) over model parameters. The model is selected, trained, and verified using Finnish Meteorological Institute (FMI) radar reflectivity composites. For the verification, both deterministic and probabilistic metrics are calculated in order to assess prediction skill against other probabilistic models and deterministic counterparts.

## 1.2 Structure of the Thesis

The present work is organized as follows. Chapter 2 contains background and a literature review on precipitation nowcasting and bayesian deep learning, aiming to familiarize the reader with essential concepts regarding the subject. Chapter 3 describes the experimental details of the work performed, including the datasets used for training and verification, the models implemented, as well as verification methods and baselines.

Chapter 4 presents the results of the experiments performed, including nowcasting examples of meteorologically interesting events, prediction skill measured with both deterministic and probabilistic metrics against multiple baseline models, as well as an assessment of predictive uncertainty. Chapter 5 discusses these results, their impact, and validity in details. Finally, Chapter 6 closes the thesis by summarizing the most important findings and takeaways.

# Chapter 2

## Background

### 2.1 Precipitation Nowcasting

#### Numerical Weather Prediction

Numerical weather prediction (NWP) is nowadays the main driver of operational weather prediction worldwide [9, 56]. On a very coarse level, NWP works by aggregating lots of meteorological observational data from sensors, which is used as initial state in solving atmospheric equations. This data can come from *in situ*, a.k.a. close contact ground-based sensors such as rain gauges and thermometers, or upper atmosphere measurements such as those from radio sondes or aircraft sensors. Additionally and most importantly, remote-sensing sources such as satellites and weather radars provide lots of observational data, and their efficient use is one of the reasons for the improvements of NWP. Then, data assimilation is performed with that data, filling in gaps and agglomerating different sources to be consistent with each others, making the output of this process appropriate to be fed as input into a numerical simulation of the atmosphere. These simulations then consist of solving Navier-Stokes equations multiple timesteps into the future. With the model output in hand, different post-processing can be applied to answer various questions about the future state of weather, such as making precipitation forecasts.

NWP can be performed on multiple spatio-temporal scales. Spatiotemporally coarse models cover a wide, sometimes global area and usually make simulations with higher *leadtimes*, meaning at further timesteps in the future. On the other hand, other models are specialized in providing finer spatiotemporal details with a smaller spatial extent. One example of this type of model is the High-Resolution Rapid Refresh (HRRR) model [4], developed by the United States National Oceanic and Atmospheric Administration (NOAA).

HRRR covers the continental United States, and has 3 kilometer spatial and 1 hour temporal resolution. Although very useful for many purposes, even such fine-detailed NWP models have not quite enough details to be most useful in very short term and very localized precipitation nowcasting as described in Chapter 1. As such, radars and radar-based nowcasting techniques shall be now described.

## Weather Radars and Products

Precipitation nowcasting is largely based on the data collected from weather radars. There are many reasons for this, the main one being that radars are capable of directly observing precipitation particles (called hydrometeors) over a significant range with a high update rate and resolution [54]. This is something that no single other ground, upper-atmosphere, or radar observation is capable of providing, especially on the mesoscale ranging from tens to hundreds of kilometers. Combining many radars into a network can further increase sensing capacity to areas covering countries such as NEXRAD covering the continental United States [2], or entire continents such as the OPERA network covering most of Europe [52].

Radar was invented as a method to perform effective military surveillance during the second world war. Their potential use as weather observation tools was discovered coincidentally when it was realized that some echos of unknown origin in these surveillance radars were indeed precipitation. After the war, weather radars started to be routinely used in affluent parts of the world to provide precipitation observations and warnings.[19]

The working of weather radars is based on emitting a very strong, directed and polarized short electromagnetic pulse, which will interact with objects and particles in the atmosphere. This interaction is scattering, which reflects back a tiny fraction of the electromagnetic energy emitted at first. A receiver then listens to these returning signals, that are often called *radar echos*. A radar typically scans the atmosphere radially at multiple azimuthal elevation angles. In addition to that angle, the vertical path of the beams, and consequently the vertical position of objects detected, are affected by the curvature of the earth and the refractive index of the atmosphere decreasing with height, bending the beam downwards. [19]

The lowest elevation angle scan is the one giving the most accurate information about precipitation actually hitting the ground, with echos further away being less reliable, as they come from targets upper in the sky. Higher elevation angle scans or products derived from a multitude of those can be used to further inform about the dynamics and development of precipitation. Combining scans and derived products from multiple radars into a

single image makes what is called a *radar composite*.

One of the most important features of weather radars is the wavelength of the signal emitted by the radar. Bigger wavelengths are less attenuated going through the atmosphere, but can only be used to detect bigger particles. Equipment using bigger wavelengths is known to be more expensive than that using smaller wavelengths. Oppositely, smaller wavelengths detect smaller particles but their operational range is limited. Weather radars are classified according to their wavelength using a band denomination. The three most common classes are S-band, C-band, and X-band radars, in the order of biggest to smallest wavelength.

In addition to meteorological echos, various non-meteorological sources are picked up by weather radars, such as solid obstacles, birds migrating, and insects. It is important to filter out or at least acknowledge those sources when preparing data for nowcasting.

*Reflectivity* ( $Z$ ) is the quantity describing the amount of signal reflected from hydrometeors in the sky to the radar receiver. It is usually described in units of decibel relative to  $Z$ , denoted  $dBZ$ . Standard weather radar beams are horizontally polarized, which produces the signals of interest when looking at radar scans for precipitation. One recent development of weather radars is the introduction of double polarization, meaning that in addition to the horizontal one, there are scans performed with vertically polarized beams. This enables estimation of precipitation type, drop shape, and helps with separating non-meteorological echos. Another important improvement to weather radars is the addition of doppler capacity, enabling estimation of target velocities from radar scans.

Converting the knowledge of (horizontal) reflectivity to that of rain rate is a highly non-trivial task and relies on approximations, because while water content is proportional to the third power of water droplet diameter, reflectivity of a single drop is proportional to its sixth power. Hence, total reflectivity depends on the drop size distribution, which depends on climatology and precipitation types. Given appropriate knowledge, the relationship between reflectivity and rain rate, commonly known as the Z-R relationship, is defined as

$$Z = AR^b \quad (2.1)$$

where  $Z$  is reflectivity,  $R$  is rainrate, and  $A$  as well as  $b$  are empirically determined parameters, usually by fitting rain gauge measurements to measured reflectivities. The most famous and widely used  $A$  and  $b$  parameters come from a study on drop size distribution from Marshal and Palmer in 1948 [35], where  $A = 200$  and  $b = 1.6$ .

### 2.1.1 Classical deterministic methods

When considering the atmosphere as a fluid, it is possible to apply tools and concepts of fluid dynamics to understanding processes happening in it, such as precipitation. In particular, one specific point of view is to consider precipitation as particles or matter moving through the forces exerted on it by the background atmospheric flow, here modeled as a liquid. This movement exerted by the surrounding medium is called *advection* and the medium dynamic itself the *advection field*. In practice, the flow consists of winds and other forces exerted on clouds and hydrometeors. Advection is described by the advection equation

$$\frac{\partial \psi}{\partial t} + \nabla \cdot (\psi \mathbf{v}) = 0 \quad (2.2)$$

Where  $\phi$  denotes the precipitation field and  $\mathbf{v}$  the advection field. With this in mind, the advection of precipitation can be observed from two distinct points of view, corresponding two different sets of coordinates. These are *Eulerian* and *Lagrangian* coordinates. The simpler one is Eulerian coordinates, which are those of the grid or space in which precipitation evolves. Lagrangian coordinates on the other hand center around particles or discrete unit being advected, so that each stays at the origin in its own set of Lagrangian coordinates, while the surroundings evolve relative to it.

A precipitation field where the only time evolution happening is the precipitation being advected has a common Lagrangian coordinate system for all precipitation units. This is never truly the case, but the idea is useful for developing nowcasting models. This concept can be formalized as that of *Lagrangian persistence*, meaning that in Lagrangian coordinates, the field stays constant. This is easy to understand by connecting the concept to its more natural counterpart of Eulerian persistence, where the precipitation field just stays constant from a motionless observer's point of view. Lagrangian persistence allows separating the time evolution of the precipitation field into an advection term and a source term, representing the divergence of the field, that is creation and weakening of precipitation. Assuming no divergence in the advection field, and using Eq. 2.2 this is formalized as

$$\frac{\partial \psi}{\partial t} + \nabla \psi \cdot \mathbf{v} = S \quad (2.3)$$

Where  $S$  denotes the source term. Basic precipitation field extrapolation based methods usually focus on modeling the advection part as accurately as possible, while methods willing to go a step further try and model the source, representing nonlinear evolution of the field.

### Extrapolation based models

The most basic form of modern model of the first category is pure extrapolation along an estimated advection field. Very often nowadays the advection field is estimated using an optical flow method, such as Lukas-Kanade optical flow [33] on successive radar frames and the precipitation field is extrapolated along that advection field. This is usually done using a Semi-Lagrangian Extrapolation scheme, which is a cheap way of performing numerical integration for the advection problem by breaking it into multiple interpolation operations.

Other models of the first category attempt to model the time evolution of the advection field and add some sort of diffusion term to better represent the loss of predictability being faster at smaller scales [24]. One such model is developed by Ryu et al. in 2020 [50]. This model models the evolution of the advection field using Burgers' equation, containing an advection and diffusion term, and additionally adds a facultative diffusion term to the nowcasted field. Another example of a similar modification is the work of Sakaino (2013) [51] where advection field temporal evolution was modeled using the Navier-Stokes with a continuity equation, and an anisotropic diffusion term, better conserving important details, was applied to the field.

Classical techniques such as Tracking Radar Echos through Correlation (TREC) [47] still being in great use. TREC simply calculates correlations between neighboring areas in successive frames and extrapolates the precipitation field along directions of maximum correlation. This technique is simple but suffers from problems in maintaining the spatial structure of echos over longer time-frames. Over the years, there has been an accumulations of models proposing improvements to TREC. There exists also many storm cell based nowcasting techniques such as TITAN by Dixon and Weiner (1993) [18]. These techniques are often also based on extrapolation and many times enable tracking and forecast of cell life-cycle features, that are particularly useful in the context of convective storms.

### Modeling non-linear evolution of precipitation fields

One of the early successful (in the sense that it presents improvement over advection-based extrapolation) attempts at modeling the time-evolution of precipitation fields is the Spectral Prognosis (S-PROG) model by Seed (2003) [57]. The model is based on the observation that precipitation feature lifetimes depend on their scale, with smaller features lasting shorter than bigger ones. This information is used by decomposing the precipitation field into a cascade of fields, each corresponding to features of a certain spatial scale.

Each of those cascade level time evolution is then modeled separately using autoregressive (AR) models. One side effect of AR models interacting is that the field gets blurred with time.

A second model taking a similar approach is ANVIL by Pulkkinen et al. (2020) [44]. Instead of single reflectivity scans or reflectivity derived products, ANVIL utilizes Vertically Integrated Liquid (VIL), an estimate of the total precipitation mass in a cloud, as a proxy to precipitation hitting the ground. ANVIL decomposes the field into a scale cascade just like S-PROG, but uses an autoregressive integrated process (ARI) to model the time-evolution at each level. ARI are autoregressive processes applied to the time derivative of the fields. This helps not to lose smaller details in the process. ANVIL was shown to surpass S-PROG for the prediction of intense precipitation ( $> 5 \text{ mm/h}$  and  $> 20 \text{ mm/h}$ ) using multiple verification metrics.

There is also examples of using phased-array weather radars to do three-dimensional extrapolation nowcasting, allowing to capture some physical processes, such as linear patterns of growth and decay. One example of this is the work of Otsuka et al. (2016) [40].

### 2.1.2 Classical probabilistic methods

Probabilistic precipitation nowcasting has high utility and therefore a multitude of models have been developed for it. Probabilistic nowcasting models can be roughly divided into two categories: Quantile-based and ensemble-based methods, the latter of which being the focus of this work. They are methods where a multiple of stochastic nowcasts are created, the set of them being called an ensemble, and probabilistic features such as rain rate exceedance probabilities are estimated from the data distribution of that ensemble. There exists also neighborhood methods, that draw new precipitation values from neighbors at random to produce a probabilistic nowcast. The first of such methods to be implemented is by Andersson and Ivarsson (1991) [6].

The most influential of modern probabilistic nowcast methods is certainly STEPS by Bowler et al. (2006) [13]. It is an ensemble method that is close to S-PROG in that it divides the precipitation field into a scale cascade and models their evolution independently through AR models. Uncertainty in STEPS is modeled through the injection of stochastic noise per ensemble member at smaller scales. One interesting feature of STEPS is that it allows blending an NWP forecast to radar images to create a more reliable nowcast at further leadtimes. STEPS is shown to retain some prediction skill up until six hours leadtime.

Another more recent ensemble nowcasting model is Lagrangian Integro-

Difference equation model with Autoregression (LINDA) by Pulkkinen et al. (2021) [43]. LINDA is designed specifically to accurately nowcast heavy localized rainfall and comes in two variants: LINDA-D for deterministic nowcasts, and LINDA-P for producing probabilistic nowcasts. LINDA is composed of multiple steps that are the identification of rain cells, advection, AR modeling for the growth and decay of features, convolutions described loss of detail at small scale, and finally stochastic perturbations in the case of LINDA-P. LINDA-D was shown to surpass both S-PROG and ANVIL in terms of skill at  $> 5 \text{ mm/h}$  and  $> 20 \text{ mm/h}$ . LINDA-P on the other hand was shown to produce more reliable and discriminative, as well as better calibrated forecasts than STEPS at high thresholds such as those described above. The forecasts are also less blurry than those of S-PROG or STEPS. In essence, LINDA describes the state-of-the-art in terms of nowcasting high-intensity localized rain.

## 2.2 Deep Learning for Nowcasting

### 2.2.1 Artificial Neural Networks

Artificial Neural Networks, abbreviated Neural Network or NN in a machine learning context, are computational systems inspired by biological neural networks, such as those present in the human brain. In this work, the terms Deep Learning will be used as meaning "regarding neural networks and their usage". Neural networks serve to accomplish many tasks, usually of the domain of artificial intelligence, such as regression, classification, or agent decision making. NN consist of discrete units called neurons linked to each others, usually arranged in functional units known as layers, in which case connections are formed between layers. Input data is fed into the NN to an input layer, transforming the input through learnable parameters, and this transformed signal is propagated through other layers further transforming the signal. The organization of layers and the way they are connected is known as the *network architecture*. Finally, signals converge to an output layer giving the product of the network, such as a prediction or class of input data. Between neurons and layers there are non-linear *activation functions*, which are in essence non-linear transformation on data, which are the basic mechanism that enables neural networks to learn complex nonlinear patterns.

## Neural Network Optimization

In Neural Networks, parameters are usually not tuned by hand due to the sheer complexity of the task. Instead, finding optimal parameters for performing the task is framed as an optimization problem. This is accomplished using standard unconstrained optimization tools and the Backpropagation algorithm, allowing learning of the parameters, despite of the problem being very high-dimensional and non-convex. This algorithm is based on producing an output by passing input data through the network, and then calculating a metric of how well the network succeeded at the task, known as a *loss function* or cost function. The gradient of the output related to network parameters are then calculated backwards starting from the loss function, flowing through the network using the chain rule of derivation. These gradients are then used to update model parameters using simple gradient descent or its variations.

Usually, the input data is divided into small subsets called *batches* or mini-batches, and the gradient updates are calculated by going iteratively through them in a process called *training*. Gradient descent over these mini-batches is called Stochastic Gradient Descent. Going through all of the training data one time is called an *epoch*, and training is usually continued for many epochs. In deep learning, available data is usually split into training, validation, and test data. Validation data is used for the validation process, in which performance of the NN is assessed on data unused for training. This information on performance can then be used to tune *hyperparameters*, that is parameters that are constant and set before training, or perform other types decision-making related to the training process, such as early stopping of the training, if performance is likely not to improve anymore. The reason why training data is not used here is that the network usually performs better on data that it was trained on compared to other data. Lastly, test data is used for independent assessment of the network performance. This is needed because even though validation data is not used for training, the network is still biased to perform better on it if decisions that improve performance were made based on it. Validation data can thus not be trusted on as an independent evaluator of network performance.

## Different Varieties of Neural Networks

There are multiple types of Neural Network architectures that are each good at different types of tasks. Main types with some relevance to nowcasting and Bayesian Deep Learning are feed-forward neural networks, Recurrent Neural Networks (RNN), and Convolutional Neural Networks (CNN). Feed-forward

neural networks are the earliest type of NN, with linear transformations of data in layers of neurons, where neurons of successive layers are all connected to each others.

RNN are a modification of feed-forward neural networks with recurrent connections along units forming a network along a temporal sequence, making them widely used for the forecasting of time-series data and natural-language processing. Important improvements over original RNN are Long-Short Term Memory (LSTM) networks, with units having internal states allowing learning longer-range temporal dependencies, and Gated Recurrent Unit (GRU) modifying LSTM units to make them more lightweight.

CNN have layers that are convolutional filters applied over the input data. CNN are great for use in computer vision and other image related tasks, as they take advantage of the fact that learning short-range dependencies is enough to perform well in many of those tasks, thus reducing the hypothesis space of learnable representations in an informed way. One CNN architecture of particular interest in the context of this work is U-Net, that has shown excellent results in e.g. semantic segmentation tasks, that are important for biomedical imaging and autonomous vehicle applications. U-Net follows a simple encoder-decoder architecture. The next section will describe current applications of artificial neural networks to precipitation nowcasting.

### 2.2.2 Deep Learning approaches to Nowcasting

There exists now many methods for producing deterministic nowcasts with machine learning. Most suffer the same problem of producing overly blurry nowcasts only after a few timesteps. Despite there having been marginal improvement with regards to this problem over the years, but it remains one important challenge when it comes to producing deterministic nowcasts. The following chapters will present some of the main advancements in Deep Learning for nowcasting since its introduction.

#### Basic Approaches

The first dedicated Deep Learning model for precipitation nowcasting is ConvLSTM by Shi et al. [59], published in 2015, building upon the work of Oh et al. [39]. ConvLSTM is a neural network combining the image feature encoding capacities of CNN and temporal prediction capacities of LSTM into a single network fusing the two. ConvLSTM was shown to outperform an optical flow based model at predicting precipitation exceeding a low rain rate of 0.5 mm/h. A further model inspired by ConvLSTM called TrajGRU is developed by Shi et al. [60] in 2017 by replacing the LSTM component with

GRU and making recurrent connections dynamically location-variant. These connections improve the predictive skill over invariant ones (ConvGRU), but the model was not tested against ConvLSTM.

Another class of models that have gained traction lately for nowcasting are simple U-Net based CNN, that adopts an image-to-image translation perspective on the problem. These networks are fed a sequence of radar images and output either one or several future frames. One example of this approach is that of Agrawal et al. [3]. In case only one frame is outputted, predictions for several leadtimes in the future may be done by iteratively feeding back predictions into the network as input, as is done with Ayzel et al.’s RainNet [8]. This iterative approach is more flexible but has the disadvantage of exacerbating the blurring problem. The CNN models in general have the advantage of being somewhat computationally less expensive when compared to ConvLSTM variants.

### Notable Recent Improvements

Recently, Pan et al. improved the nowcasting of convective evolution by adding polarimetric variables (ZDR and KDP) in addition to reflectivity scans as inputs to a U-Net based model [41]. These polarimetric variables are derived from double polarization weather radars and their values are strongly associated with life-cycles of convective cells. This work showcases the importance of multichannel data and not only relying on reflectivity if one wants to accurately model nonlinear evolution of precipitation.

Above models suffers badly the the problem of blurring stated above, as do all current models based on discriminative neural networks, which are networks mapping one input to one output. Prediction blurring is in part related to loss functions used, as losses like  $\ell_2$  (Mean Squared Error (MSE)) and  $\ell_1$  tend to act that way when minimized with some uncertainty present. This excessive blurring also hinders the prediction of useful patterns such as heavy localized rainfall. Multi-Scale Structural Similarity Index (MS-SSIM), originally an image quality assessment metric [65], has recently started to be used as a loss function in deep learning, particularly in image reconstruction tasks. One of its main assets is that its single-scale counterpart (SSIM) has been shown to reduce blurring in image reconstruction [68] making it a good candidate loss function for nowcasting with neural networks. Indeed in recent years, there has been some cases where MS-SSIM or SSIM started have started to be used in Deep learning based nowcasting models, such as the one of Yin et al. (2021) [67]. Here both SSIM and MS-SSIM produced results surpassing MSE, with MS-SSIM giving the best results of the two. Most importantly, these loss functions seemed to better preserve high rain

rates, which tended to vanish with traditional losses at higher leadtimes.

Perhaps one of the biggest breakthroughs made yet, is the use of Generative Adversarial Networks (GAN) by Ravuri et al. last year [46]. This approach of using a generative model, i.e. one that generates samples from a probability distribution conditioned on past radar measurements, manages to solve the problem of nowcast blurring. GAN are a class of networks based on the competition between a generator network that generates the samples and discriminator(s) networks trying to discern whether these generated samples are real or fake. The generator then tries to fool the discriminators in a zero-sum game. Ravuri et al. use a generator network based on a convGRU architecture, nowcasting 90 minutes at once, two discriminators, and a regularization term penalizing deviations of generated samples at the grid level. The first discriminator ensures spatial consistency and the lack of blurring, while the second one ensures temporal consistency in generated sequences. The resulting generative model has slightly superior skill compared to existing approaches, but most importantly overwhelmingly over alternatives preserves features that make for a good nowcast from an operational forecaster point of view and is capable of providing probabilistic nowcasts which is very useful [46]

There has so far been limited work in probabilistic nowcasting with Deep Learning. In addition to the GAN by Ravuri et al. that can be used for such purpose, a probabilistic nowcasting model called MetNet was developed by Sonderby et al. [62]. The model is based on Axial Self-Attention by Ho et al. [27] and is capable of beating HRRR on probabilistic verification metrics for leadtimes of up to 8 hours.

## 2.3 Bayesian Deep Learning

This section gives an overview on Bayesian Neural Networks, that will be used make probabilistic nowcasts in this work.

### 2.3.1 Learning Probability Distributions

Neural networks are extremely useful model that on the other hand are very complex, in the sense that they have many optimizable parameters. The effect of this is that they are sensible to *overfitting*, meaning that if left unchecked they will learn spurious patterns in the data, over-adapting to the training dataset and worsening generalization ability. In order to counter this, different mechanisms exist that bias the neural network towards learning simpler representations that have better generalization ability when pre-

sented with out-of-training-data examples. This concept is known as *regularization*. Over the past decades, many regularization mechanisms have been successfully developed and applied to the training of deep neural networks. Notable examples include Dropout and weight decay, also known as L2-regularization.

One caveat of these classical regularization methods is that they do not enable representing the uncertainty of neural networks in unexplored regions, although they do improve performance in them. As it has understandably many benefits to make a neural network able to say "I don't know", a way to represent different plausible scenarios arising with novel data is needed.

There are two ways of approaching uncertainty estimation in neural networks. One is to directly model the uncertainty of predictions, and the other is to model the uncertainty of model parameters. These two are often complementary, but The latter approach has the benefit of providing implicit regularization to the network and is indeed the primary focus of this work. Learning this uncertainty of model parameters is most often accomplished using Bayesian Neural Networks (BNN). Strictly speaking, BNN are a class of stochastic neural networks, that is NN with stochastic components, where the parameters are probability distributions that are estimated using Bayesian inference. Bayesian Neural Networks in their current form were introduced by MacKay in 1992 [34], after early works starting in 1987 by Denker et al. [17], Tishby et al. [63], Denker and LeCun [16], and Buntine and Weigend [14].

In a Bayesian framework, the posterior distributions of parameters are estimated conditioned on data and a prior. Because exact Bayesian inference involves dealing with intractable integrals, it is not doable to solve them for real-world neural networks having thousands to millions of parameters. As such, two broad categories of inference methods have been developed for use in BNN. The first one is using Markov-Chain Monte-Carlo (MCMC) to sample from posterior distributions. MCMC works well for smaller-scale models, but it is limited to only thousands of parameters because of the computational expensiveness of Markov-Chain simulations.

The second inference method category is *Variational Inference* (VI). The main idea behind variational inference is to turn the problem of finding the intractable true bayesian posterior into that of finding the closest distribution from a limited distribution family (the variational family) with regards to the true posterior. This turns the problem of solving an integral into that of optimization, allowing the use of classical unconstrained optimization methods.

Dropout was surprisingly shown to be an instance of variational inference with a Bernoulli distributed posterior by Gal and Ghahramani in 2016 [23]. This legitimizes the use of Monte-Carlo Dropout, a technique used to get predictive uncertainty estimates from networks with Dropout layers, by simply

not switching off the layers for inference and thus drawing Monte-Carlo samples of the data distribution. Although not as expressive as explicit Bayesian NN, MC dropout has no computational overhead and is simple to implement, making it a very popular alternative to BNN.

In the context of Bayesian Deep Learning, the underlying neural network architecture will be referred to as the functional architecture or the functional model. Other components such as the inference method, the prior, or the posterior modeling will be commonly referred to as the stochastic model.

### 2.3.2 Variational inference

Variational inference as a means of training bayesian neural networks was first introduced by Hinton and Camp in 1993 [26]. However it was not used for a long time before breakthroughs by Graves et al. (2011) [25] and Blundell et al. (2015) [12] permitted its use in large-scale neural networks. The following sections will introduce the loss function and optimization algorithm used, mostly based on the work of Blundell et al. [12].

#### Evidence Lower Bound loss function

Finding the variational posterior  $q(\mathbf{w}|\theta)$  best approximating the true posterior is formalized as minimizing the Kullback-Leiber (KL)-divergence between the two distributions as

$$\theta^* = \arg \min_{\theta} \text{KL}(q(\mathbf{w}|\theta) || P(\mathbf{w}|\mathcal{D})) \quad (2.4)$$

where  $\theta$  refers to variational posterior parameters,  $\mathbf{w}$  to the weights,  $\mathcal{D}$  to the data,  $q$  to the variational distribution, and  $P$  to the true distribution. The KL-divergence is defined as

$$\text{KL}((q(\mathbf{w}|\theta) || P(\mathbf{w}|\mathcal{D}))) = \int q(\mathbf{w}|\theta) \log \frac{q(\mathbf{w}|\theta)}{P(\mathbf{w}|\mathcal{D})} d\theta \quad (2.5)$$

This definition can be used to turn the minimization objective into

$$\arg \min_{\theta} \mathbb{E}_{q(\mathbf{w}|\theta)} [\log q(\mathbf{w}|\theta)] - \mathbb{E}_{q(\mathbf{w}|\theta)} [P(\mathbf{w}, \mathcal{D})] + \log P(\mathcal{D}) \quad (2.6)$$

Here the evidence  $P(\mathcal{D})$  is very difficult to estimate. Luckily though, the minimization objective does not depend on it, so a new objective called the evidence lower bound (ELBO) or variational free energy is used instead, which is defined as

$$\text{ELBO}(q) = -( \text{KL}((q(\mathbf{w}|\theta) || P(\mathbf{w}|\mathcal{D}))) - \log P(\mathcal{D})) \quad (2.7)$$

Solving for the evidence  $\log P(\mathcal{D})$ , it is clear that ELBO is a veritable lower bound for the it, as KL-divergences can not be negative.

ELBO can be rewritten as

$$\mathcal{F}(\mathcal{D}, \theta) = \text{KL}[q(\mathbf{w}|\theta)||P(\mathbf{w})] - \mathbb{E}_{q(\mathbf{w}|\theta)}[\log(P(\mathcal{D}|\mathbf{w}))] \quad (2.8)$$

which is a form more suitable for the development of an optimization method. Here and from now on, ELBO will be referred to with the  $\mathcal{F}$  symbol. For the rest of this work, the variational posterior family will be assumed to be that of Gaussian distributions, as the *Bayes-by-Backprop* (BBB) algorithm described in the following section works on this posterior family.

### The Bayes-by-Backprop Algorithm

In order to apply backpropagation to Bayesian Neural Networks with variational inference, any sampling operation must be made external to the network in order to allow gradients to flow backwards. This is accomplished by parametrizing the standard deviation  $\sigma$ , which is the parameter linked to uncertainty in weights, with a parameter  $\rho \in \mathbb{R}$ , such that  $\sigma = \log(1 + \exp(\rho)) \in (0, \infty)$ . This is known as the reparametrization trick. Weights  $\mathbf{w} = t(\sigma, \epsilon)$  are made deterministic functions of posterior parameters and the external stochastic noise  $\epsilon$ . This is what allows gradients to flow through the network by making sampling external.

Using the reparametrization trick and the proposition

$$\frac{\partial}{\partial \theta} \mathbb{E}_{q(\mathbf{w}|\theta)}[f(\mathbf{w}, \theta)] = \mathbb{E}_{q(\epsilon)}\left[\frac{\partial f(\mathbf{w}, \theta)}{\partial \mathbf{w}} + \frac{\partial \mathbf{w}}{\partial \theta} \frac{\partial f(\mathbf{w}, \theta)}{\partial \mathbf{w}}\right] \quad (2.9)$$

from Blundell et al. [12], the ELBO cost function (Eq. 2.8) can be rewritten as in a form amendable to minibatch optimization, yielding

$$\mathcal{F}(\mathcal{D}, \theta) \approx \sum_{i=1}^n \log q(\mathbf{w}^{(i)}|\theta) - \log P(\mathbf{w}^{(i)}) - \log P(\mathcal{D}|\mathbf{w}^{(i)}) \quad (2.10)$$

an unbiased Monte-Carlo estimator of ELBO. Here  $\mathbf{w}^{(i)}$  denotes the  $i$ :th monte-carlo sample drawn the the posterior. The proposition 2.9 follows from the reparametrization trick.

From the terms in eq. 2.10:

1.  $\log q(\mathbf{w}^{(i)}|\theta)$  is the log likelihood of weights given the current posterior distribution.

2.  $-\log P(\mathbf{w}^{(i)})$  is the negative log likelihood of the weights given the prior, which is usually not learned and serves as a regularization mechanism.
3.  $-\log P(\mathcal{D}|\mathbf{w}^{(i)})$  is the likelihood term, dependent on the data  $\mathcal{D}$

The first two terms are grouped together as the complexity cost, while the last term is often called the Likelihood cost. With all of the conditions in place, optimization takes place in pretty much the same way as for basic backpropagation, only updating two variables instead of one ( $\mu$  and  $\rho$  instead of point estimates). One step of the optimization loop is

1. Sample  $\epsilon \sim \mathcal{N}(0, 1)$
2. Let  $\mathbf{w} = \mu + \log(1 + \exp(\rho)) \circ \epsilon$ ,  $\theta = (\mu, \rho)$   
◦ referring to point-wise multiplication.
3. Let  $f(\mathbf{w}, \theta) = \log(q(\mathbf{w}|\rho)) - \log(P(\mathbf{w})P(\mathcal{D}|\mathbf{w}))$
4. Calculate gradients w.r.t.  $\mu$  and  $\sigma$ :
  - (a)  $\Delta_\mu = \frac{\partial f(\mathbf{w}, \theta)}{\partial \mathbf{w}} + \frac{\partial f(\mathbf{w}, \theta)}{\partial \mu}$
  - (b)  $\Delta_\rho = \frac{\partial f(\mathbf{w}, \theta)}{\partial \mathbf{w}} \frac{\epsilon}{1 + \exp(-\rho)} + \frac{\partial f(\mathbf{w}, \theta)}{\partial \rho}$

Then parameters are simply updated as:

- (a)  $\mu \leftarrow \mu - \alpha \Delta_\mu$
- (b)  $\rho \leftarrow \rho - \alpha \Delta_\rho$

### Complexity cost weighting and priors

The cost to be minimizing can be again rewritten as

$$\mathcal{F}_i^\pi(\mathcal{D}, \theta) = \pi_i \text{KL}[q(\mathbf{w}|\theta) || P(\mathbf{w})] - \mathbb{E}_{q(\mathbf{w}|\theta)}[\log P(\mathcal{D}_i|\mathbf{w})] \quad (2.11)$$

Here  $\pi_i$  is the relative weighting of the complexity cost in the  $i$ :th mini-batch. There are many ways to weigh the complexity cost against the likelihood cost, but an usual constraint is that  $\pi \in [0, 1]^M$  and  $\sum_{i=1}^M \pi_i = 1$ . Graves et al. (2011) [25] used equal weighting as  $\pi_i = 1/M$ , but Blundell et al. (2015) [12] found the scheme  $\pi_i = \frac{2^{M-i}}{2^M - 1}$  to offer better performance in their experiments.

The prior distribution  $P(\mathbf{w})$  is often chosen to be a diagonal Gaussian distribution, because it allows calculating the KL-divergence with regards to the Gaussian posterior analytically. A Gaussian prior placed on weights

would be equivalent to L2-regularization, also known as weight decay. Recently, despite of the inability to analytically calculate KL-divergences using it, Gaussian scale mixture priors have also started to be used [12, 61] because of their properties facilitating optimization. These scale mixture priors, if containing two scales are defined as

$$P(\mathbf{w}) = \prod_j \alpha \mathcal{N}(\mathbf{w}_j | 0, \sigma_1^2) + (1 - \alpha) \mathcal{N}(\mathbf{w}_j | 0, \sigma_2^2) \quad (2.12)$$

where  $\sigma_1 > \sigma_2$  and often  $\sigma_2 \ll 1$ ,  $\alpha$  is the relative weights of the two scales, and  $\mathbf{w}_j$  is the  $j$ :th weight of the neural network.

### Alternative reparametrizations

One problem with conventional Bayes-By-Backprop is that the same  $\epsilon$  is used for each parameter in a layer. One side effect of this is that gradients of different weights are correlated, hindering the training.

Kingma et al. (2015) [30] introduced a modification to the reparametrization trick, called the *local reparametrization trick* (LRT). This modification works similarly but rather than weights, pre-activation layer outputs are sampled using  $\epsilon$ , with a different  $\epsilon$  for each output. This reduces variances in a minibatch, allowing for faster training overall. Kingma et al. showed that the unbiased Monte-Carlo estimator of log likelihood in Blundell et al. has variance that does not decrease with batch size because of the contribution of batch member covariances due to shared  $\epsilon$ , and consequently designs an estimator with zero covariance, leading to the method above. Gradient variances with LRT thus are inversely proportional to batch size, leading to easier optimization with higher batch sizes.

Flipout by Wen et al. [66] is another modification to Bayes-by-Backprop, attempting to solve the same problem. While LRT can only be used for fully-connected feed-forward neural networks with no weight sharing, Flipout can be used on other architectures too [66]. Flipout works the same as BBB, except for the fact that it multiplies the sampled  $\epsilon$  by a random sign matrix of the size of the parameter space. This way the  $\epsilon$  are to some degree made pseudo-random, decreasing the gradient variances for a very meager computational cost.

### 2.3.3 Predictive uncertainty estimation and decomposition

*Predictive uncertainty* is defined as the total uncertainty arising in practice when making predictions with a probabilistic or ensemble based model.

Broadly speaking, in deep learning predictive uncertainty can be divided into two main categories. These are *aleatoric uncertainty*, which originates from the input data, and *epistemic uncertainty*, which originates in the inaccuracy of the model. Epistemic uncertainty can be reduced with more training or some other changes, while aleatoric uncertainty can not be by definition. Furthermore, aleatoric uncertainty is classified into homoscedastic and heteroscedastic aleatoric uncertainties in the context where one tries to estimate it. The former assumes that all observation share the same underlying uncertainty, while the latter allows each observation uncertainty to differ. [61]

Characterization of uncertainty is important because not taking it into account might lead to over-estimation or under-estimation of failure probability of a model, as explained by Der Kiureghian [31]. There exists techniques to decompose predictive uncertainty into aleatoric and epistemic uncertainties in deep learning. To model the epistemic component, these techniques base themselves on ensembles just like those produced by Bayesian Neural Networks, or the variability between ensemble members to be more precise. Strategies for modeling the aleatoric part differs between classification and regression tasks. In classification, heteroscedastic aleatoric uncertainty can be directly inferred from class logits without any network modification [32, 61]. For regression on the other hand, there is a need to encode the uncertainty in the data as it is not there in the first place. By modeling the likelihood as homoscedastic Gaussian likelihood, one can learn a common homoscedastic uncertainty term for the dataset. By modeling the data as heteroscedastic Gaussian likelihood on the other hand, one can learn a different uncertainty term for each data point. This is done by separating the output of the network and its aleatoric uncertainty into two different channels towards the end of the network, and plugging these in the Gaussian likelihood accordingly [29].

### Uncertainty in Radar-based Nowcasting

In radar-based nowcasting, the aleatoric uncertainty of models can be divided into that emanating from measurements and that coming from other factors induced by the processing of data inside the nowcasting process. Measurement uncertainty can again be subdivided into sampling-based and system-based uncertainties. Sampling uncertainty refers to the randomness coming from which hydrometeor or obstacle the radar beam encountered in the scanning process, and system-related ones come from various factors in the radar itself. According to Cao et al. [15], sampling-based uncertainties dominate in well-configured radar systems. Other factors inducing aleatoric uncertainty include all stages where information is possibly lost, such as data compression

operations. Additionally, for methods where model outputs are iteratively fed back into the algorithm to produce new predictions, the whole compound uncertainty of the previous step output is accounted as aleatoric uncertainty in the next step.

Prediction error and uncertainty are related to each others but not equivalent. Sources of errors can be used to partly infer sources of uncertainty, keeping in mind that model outputs may have high bias (possible error) but low variance (uncertainty). Bowler et al. enumerates sources of prediction errors in the context of advection based models, dividing them into three categories. These are errors in estimating the initial advection field, in modeling its time evolution and the Lagrangian evolution of features [13]. Part of these errors are related to the model itself, and as epistemic uncertainty can by definition be reduced, these sources of error can be decreased in practice. This can be achieved by for example improving data quantity and quality with regards to the task at hand, improving the functional model to better express dependencies between outputs and outputs, and by improving the training procedure, such as a more fitting prior in Bayesian Deep Learning.

# Chapter 3

# Materials and Methods

From now on are presented the data and models used for experiments, and how those nowcasting models were verified.

## 3.1 Experimental Setup

In practice, predictions are calculated for BCNN, described in section 3.3.2, as well as all models of sections 3.4.1 and 3.4.2 using all the timestamps contained in the test set described in section 3.2. Predictions are computed for 36 timesteps, which is equivalent to 3 hours into the future with an interval of 5 minutes. After that, deterministic verification metrics described in section 3.5.1 are computed for using the predictions of each model, and probabilistic metrics from section 3.5.3 are calculated from BCNN and probabilistic models of section 3.4.2. Then metrics are averaged over the set for each leadtime of interest.

In order to make sure that results are valid, all timestamps having any (even a single) observation missing are discarded, and similarly timestamps having any prediction from any model missing are also removed from calculations. Additionally, only pixels where data is present in the predictions of all models are counted in metric calculations. This is accomplished by calculating a common NaN mask using the logical OR operation over NaN values of each model, and subsequently applying that common mask to each prediction.

Regarding probabilistic forecasts, uncertainties present in the data are coming from a diverse set of sources, including both aleatoric and epistemic uncertainties. In order to accurately represent the data distribution resulting from this compound uncertainty, a large ensemble size is needed. Hence, the ensemble size was chosen to be 48 for final models. Verification might be even

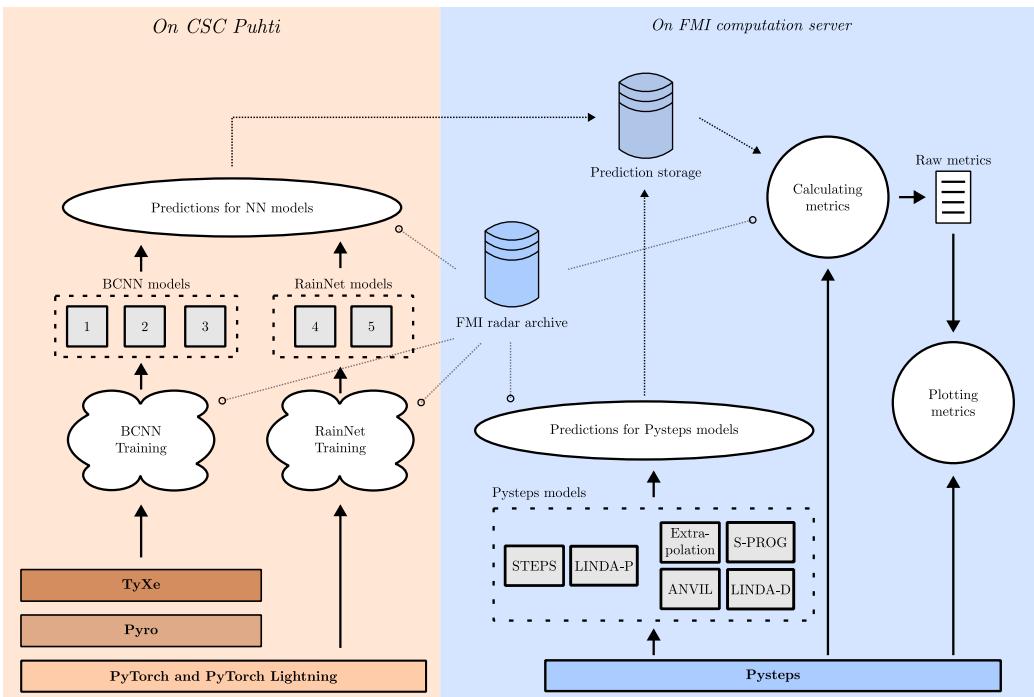


Figure 3.1: Overall workflow used for this work, tied with used libraries for different components.

more reliable with bigger ensembles, but this would be at the expense of too much storage space needed for predictions, which would be very inconvenient. Chosen sizes were deemed to be a good compromise regarding this.

Intermediate results are saved such that predictions are saved in HDF5 archives, in a format where each predicted radar image is a separate dataset in deterministic models, and each ensemble of predicted radar images for a certain leadtime is its separate dataset in ensemble nowcast models. In this storage procedure, predicted float reflectivity values are first compressed into UINT8 with a scale-offset scheme, then thresholded to 0.1 mm/h to facilitate further lossless compression which is finally carried out using GZIP. Trying to reduce the storage space needed is essential because predictions take a very large amount of space on disk. This is a consequence of combining a big number of timestamps, ensemble members and leadtimes for multiple models. Raw metrics do not suffer from the same problem and are saved to disk in a binary Numpy format.

### 3.1.1 Software and resources

The deep learning models were all implemented using the PyTorch framework and the PyTorch Lightning wrapper [20]. These libraries were chosen because of the combined ease of prototyping brought by PyTorch Lightning, and the maturity and flexibility of the parent framework PyTorch. RainNet was ported to PyTorch following the original Tensorflow implementation by Ayzel et al. [7].

As for the implementation of the probabilistic inference mechanisms for Bayesian Neural Networks, choice was made not to implement them by hand, but to rely on the machinery contained in the Probabilistic Programming Language (PPL) Pyro [10], which is itself built on top of PyTorch and includes a fully-featured implementation of Stochastical Variational Inference (SVI). In order to facilitate the implementation task, the TyXe package [48] was used. TyXe is a library designed to provide an interface simplifying the implementation of Bayesian Neural Networks using PyTorch and Pyro. A few of the reasons why TyXe was chosen are that it permits easily turning existing neural networks into BNNs without having to use hard-coded bayesian layers, and dynamically switching on-and-off the local reparametrization trick and flipout in layers. Some problems encountered include components of TyXe having difficulties working together with PyTorch Lightning abstractions.

Verification experiments and non-deep-learning baseline models were ran and implemented using the open-source library Pysteps [45]. It provides implementations for all non-deep-learning models described as well as implementation of verification metric primitives used in this work, and tools for their visualization.

The computational resources from the Finnish IT Center for Science (CSC) were used for GPU intensive tasks such as training and calculating predictions with deep-learning models, and one of FMI’s computational servers was used for performing other, more CPU-intensive operations such as predicting with baseline non-deep-learning models and calculating verification metrics. We trained the models on the CSC Puhti supercomputer, using one Nvidia V100 GPU with 32GB of VRAM, 64GB of RAM, and 10 cores from a 2.1 GHz Intel Xeon Gold 6230 CPU. As for the FMI server, it contains 2 Intel Xeon Gold 6138 2.0 GHz CPUs with each 20 cores and 2 threads by core, with 192GB of RAM.

## 3.2 Dataset and data selection

As input data we use lowest elevation angle radar reflectivity composites with 1km spatial resolution and 5 minute temporal resolution from the Finnish Meteorological Institute, cropped into a 512x512 km region covering southern Finland. The domain and bounding box are illustrated in Figure 3.2. The radar network consists of C-band dual polarization radars.

Because lowest elevation angle horizontal reflectivity is a good estimator of precipitation at ground level, it is a natural data choice for building a nowcasting model. Radar composites archived from these radar scans are readily available at FMI which facilitated their retrieval for this work. Although the bounding box covering southern Finland did have some data quality issues, it was still deemed to be the best choice, when compared to other candidate datasets such as TAASRAD19 [22], as the problems were minor and did not disturb the training process. For example, missing pixels were correctly labeled, which allowed trivial removal of composites with some included.

Some of the problems with the chosen data are related to lackluster data quality control. Specifically, dual polarization information was not used for the removal of non-meteorological echos, making this removal less accurate. Also, again related to the lack of vertical polarization information, attenuation correction was not performed, meaning that echos originating from behind other strong echos will be attenuated and show as weaker than they really are. This is partly mitigated by having multiple radars in the composite, but even this will not cover all needed scenarios, making this lack of attenuation correction a significant problem.

The dataset was chosen so that at first, a selection of rainy days were chosen as to correspond to the 100 days with the most pixels exceeding a 35 dBZ reflectivity threshold during the summer period spanning from May to September during years 2019, 2020, and 2021. The present threshold was chosen as a value which convective storms usually exceed during their whole lifetime [64], as the events that are of the most interest in the context of this work are such convective storms.

This dataset is then cleaned and filtered, after which it is divided into training, validation, and test sets. Cleaning the data involves first going through all timestamps and removing those with either partially or completely missing data. The filtering part consists of removing timestamps with less than 1% of pixels containing reflectivity values exceeding 20 dBZ. Splitting of data into training, validation, and testing sets is performed by using a block sampling strategy [56] with 6 hour long blocks to prevent auto-correlation between consecutive radar images from invalidating independence

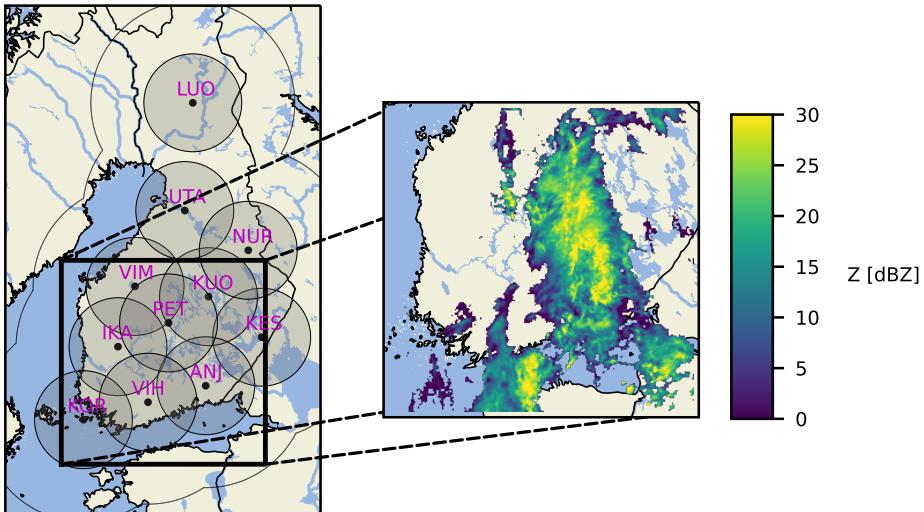


Figure 3.2: FMI radar domain and the chosen bounding box covering southern Finland. The three letter codes identify the radars, inner circles their worst case effective range (usu. during summer), and outer circles their best case effective range (usu. during winter). The lowest level reflectivity composite from the 30th of June 2020 at 00:05:00 is depicted on the right.

between splits. The final split sizes were 15840 radar images for the training split, 2664 for the validation split, and 2448 for the test split.

From radar reflectivities, rainrate estimates were calculated using Eq. 2.1 solved for  $R$ , with empirically determined parameters  $A = 223$  and  $b = 1.53$ , and  $z = 10^{Z/10}$ , giving us a relation of

$$R = (10^{Z/10}/223)^{1/1.53} \quad (3.1)$$

for the current data.

### 3.3 Model

#### 3.3.1 The baseline: A U-Net based CNN

- why not convLSTM ??? often used too!

For the implementation of a Bayesian Convolutional Network, we use as our base functional model the same architecture as for RainNet by Ayzel et al. [8], which is itself heavily inspired by U-Net and SegNet model families.

RainNet just like those families adopts a U-shaped encoder-decoder architecture. The encoder is composed of successive pooling and convolutional layers, reducing image sizes passed to the next layer while increasing the number of filters. The decoder part adopts a mirror archictecture of successive upsampling and convolutional layers, increasing the image size while reducing the number of filters. There is 5 levels in the encoder-decoder structure, just like in SegNet. Additionally, there are skip connections from the encoder to the decoder branch, carrying higher-level filter maps through, just like in U-Net. This is done in order not to lose smaller spatial scale details with pooling.

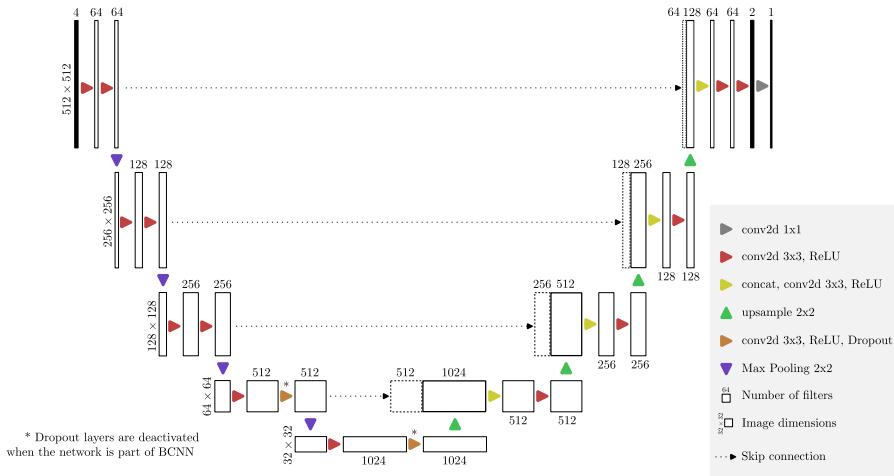


Figure 3.3: The functional CNN architecture used for this work, Identical to that of RainNet from Ayzel et al. [8]. **concat** refers to the concatenation of feature maps along the channel dimension.

The convolutional filters have a filter size of 3, with padding designed to preserve image shape, and a stride of one. Finally, the last convolutional layer in the network has a filter size of 1 with no padding, and it is followed by a linear activation output layer, i.e. no (nonlinear) activation function. The network does not contain any fully-connected layers, and in total consists of 31.4M parameters. RainNet works by feeding in 4 consecutive radar images, with the physically speaking temporal dimension represented in the channel dimension of the tensor. As the output of the network, we get one output radar image, representing the next predicted frame. In order to make predictions at further leadtimes, the predictions are fed back to the inputs one-by-one in an iterative process.

The network is trained by calculating a loss function between observed and predicted radar images. This loss function is originally the LogCosh loss, as described in Ayzel et al.’s paper [8]. In addition to that loss, the use of Multi-Scale Structural Similarity Index (MS-SSIM) [65] as a loss function was also implemented in this work. MS-SSIM is a loss function, originally an

The rainrates were scaled using a logarithmic transformation  $x_{\log} = \ln(x + 0.01)$  for use in the CNN and BCNN models, and they were additionally shifted upwards by a factor of 5 and then scaled down by a factor of 10 to fit into the range of 0 to 1 when using MS-SSIM loss. One benefit of using MS-SSIM instead of Logcosh loss is that MS-SSIM better preserves structural details of predictions such as edges and gradients, which helps with nowcasting higher reflectivity values. This same behaviour however often produces diverging nowcasts, characterized by loss of total precipitation mass and unphysically high precipitation intensities at further leadtimes. The logarithmic transformation of the data is needed because rainrates in themselves are lognormally distributed. This transformation makes them normally distributed, and thus more suitable inputs to neural networks.

The diverging behavior is alleviated by calculating the loss function for nowcasts done over several leadtimes, making the resulting learned network more temporally stable. Hence for training, the predictions are calculated for 6 leadtimes (30 minutes ahead) and the loss is calculated such that each prediction leadtime is weighted equally. Using multiple leadtimes is also implemented for the LogCosh loss function. MS-SSIM was calculated using a data range parameter of 1.0 and equal scale weights of  $[0.2, 0.2, 0.2, 0.2, 0.2]$  ordered from smallest to biggest, as opposite to the scaling optimized for visual perception introduced by Wang et al. [65] in the context of using MS-SSIM as a visual quality assessment tool.

The reason for adopting RainNet as the functional architecture of the Bayesian Neural Network is that compared to for example ConvLSTM based models, it has faster inference without losing much in skill. This is important when building upon a model, especially when having limited resources, as added components make the training and inference slower.

### 3.3.2 BCNN: a Bayesian extension to RainNet

**!FIXME BCNN is a working name, find potentially better name**  
**!FIXME!**

BCNN is the Bayesian extension made for this work of the U-Net based architecture described in section 3.3.1. In the network, posterior probability distributions of weights and biases are modeled as diagonal Gaussian distributions and optimization is carried out using variants of the Bayes-by-

Backprop algorithm described by Blundell et al. (2015) [12], minimizing the ELBO loss function as described in section 2.3.2. Posteriors for all parameters are initialized identically to a mean  $\mu$  of 0 and a variance  $\sigma$  of 1e-4. Parametrizing the posteriors as Gaussian distributions doubles the number of learnable parameters to 62.8M. The Local Reparametrization Trick (LRT) [30] and Flipout [66] are both used in an attempt to reduce gradient variance and speed-up optimization.

As for the prior distribution, two variants are implemented. One being a diagonal Gaussian prior  $P(\mathbf{w}) \sim \mathcal{N}(0, 0.1)$  and the second a Gaussian scale mixture prior as defined in Eq. 2.12, with  $\alpha = 0.5$  and  $\sigma_1, \sigma_2 = 0.3, 0.03$ . The scale-mixture prior might enable faster initialization and more accurate asymptotic behavior, but the Gaussian prior has more convenient mathematical properties, as it enables calculating KL-divergences in closed form when combined with a Gaussian posterior as in here, using the analytical expression

$$D_{KL}(q(\mathcal{D}|\mathbf{w})||P(\mathbf{w})) = \frac{1}{2}[\log \frac{|\sigma_P|}{|\sigma_q|} - d + (\mu_q - \mu_P)^T \sigma_P^{-1}(\mu_q - \mu_P) + \text{tr}\{\sigma_P^{-1} \sigma_q\}] \quad (3.2)$$

where  $q$  and  $P$  denote the multivariate diagonal posterior and prior,  $\mu_q$  and  $\mu_P$  their respective means,  $\sigma_P$  and  $\sigma_q$  their respective standard deviations and  $d$  the identity matrix of the number of dimensions equal to that of the distributions.

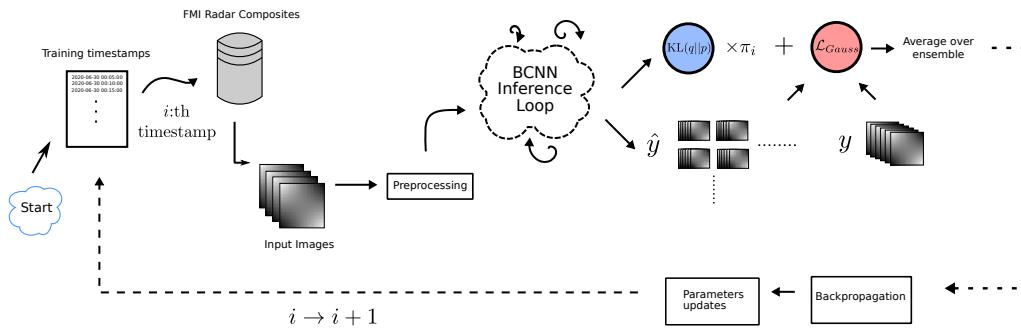
In BCNN, the radar images were scaled between 0 and 1 using the same procedure as described for the MS-SSIM loss function in the context of RainNet in section 3.3.1. The data likelihood cost was modeled as Homoscedastic Gaussian likelihood [29], defined by

$$\mathcal{L}_{Gauss}(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2\sigma^2} ||y_i - \hat{y}_i||^2 + \frac{1}{2} \log \sigma^2 \quad (3.3)$$

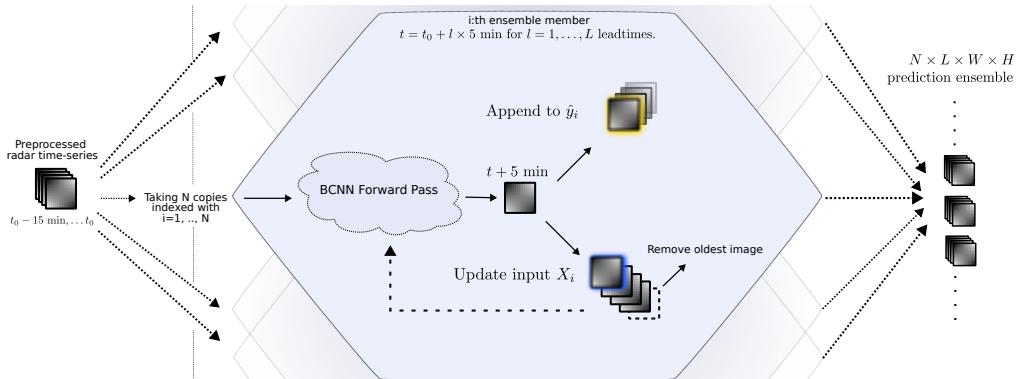
where  $i = 1 \dots N$  indexes the number of pixels in the images, and  $\sigma$  here refers to the observational noise parameter. In homoscedastic gaussian likelihood, aleatoric uncertainty, characterized by  $\sigma$  is assumed to be the same for all the data. This parameter can be learned or left constant. The latter is done in the present work, and a value of 0.02 is chosen as a guess, aiming to represent the uncertainty inherent to radar images.

Multiple procedures for weighting the complexity cost against the likelihood cost, some presented in Section 2.3.2, are implemented. In addition to the equal weighting scheme from Graves et al. [25] and the batch index

dependent scheme from Blundell et al. [12], a scheme reducing the weight of the complexity cost each epoch, defined as  $\pi_j = 2^{-j}$ , where  $j = 0, 1, \dots$  is the current epoch, is implemented to provide better potential convergence, as the equal weight scheme sometimes provides excessive regularization and the second scheme has problems with converging when using smaller batch sizes as in this work.



(a) Training procedure for the BCNN illustrated. **Rough first version: cleaning up later with radar imgs/preds over black boxes, alignment, sizes, fonts same everywhere after getting feedback on overall structure**



(b) Inference procedure for generating ensemble nowcasts with the BCNN. **Same comments as above...**

Figure 3.4: Training and inference procedures

The predictions are made by simply performing a forward pass through the network, which samples the parameters and produces a Monte-Carlo (MC) sample, i.e. an ensemble member. Predictions for further leadtimes are made in the same way as with the deterministic RainNet, that is by iteratively re-plugging the output, now a MC sample, into the inputs of

the network. It is important to note that stochasticity is preserved because parameters are re-sampled at each step of the iterative prediction. Ensembles are produced by repeating this process multiple times independently, that is, such that one initial MC sample will lead to no more than a single new MC sample each leadtime. Sampling is done this way to avoid the problem of exponential growth that would appear if each sample at leadtime  $t$  would yield more than one sample at leadtime  $t + 1$ . The inference procedure is illustrated in Figure 3.4b.

Training is performed such as to train the model to predict on a single leadtime (+5min), and after this performing further training on 6 leadtimes (+30min) until model convergence. This is done in an attempt to speed up model convergence compared to training the model on many leadtimes right from the start. At each training step, one MC sample is drawn from the network and gradients are calculated from its ELBO loss. A diagram illustrating the training process is shown in Figure 3.4a.

The model selection process is performed hierarchically as follows: Two priors per category, that is a Gaussian prior with  $\sigma = 0.1$  and  $\sigma = 0.001$ , and a Gaussian Scale Mixture (GSM) prior with  $\sigma_1, \sigma_2 = 0.3, 0.03$  and  $\sigma_1, \sigma_2 = 0.1, 0.001$  are selected and models are trained with them. Next, the best performing model is selected based on the validation likelihood cost, and it is retrained using the KL-annealing strategy proposed as well as that from Blundell et al. Finally, the best performing model overall is selected for verification against baselines.

### 3.3.3 Additional Hyper-parameters for NN

Both for RainNet and BCNN, the Adam optimizer is used with an initial learning rate of 1e-04, without any learning rate scheduler. For RainNet, other hyper-parameters were left to their `torch.optim.Adam` default values, while for BCNN, the momentum parameter  $\beta_1$  was increased to 0.95 to better accommodate the increased stochasticity of SVI compared to non-Bayesian NN [1]. Early stopping was set with condition that the validation loss does not decrease for 3 epochs for both models. For BCNN, only the likelihood part of the loss was taken into account in the early stopping criterion.

## 3.4 Verification Baseline Models

### 3.4.1 Baseline Deterministic Models

In order to verify the skill of BCNN, multiple deterministic models were used to produce nowcasts against which those produced by BCNN are compared. First and foremost, the ensemble averages of BCNN are compared against those of RainNet trained with 6 leadtimes, and both LogCosh and MS-SSIM loss, with the configuration outlined in sections 3.3.1 and 3.3.3. Although loss functions used for training the network are not comparable, it still is interesting to compare the skill between the classical predictions and Bayesian version ensemble means, as underlying functional architectures are the same.

In addition, predictions were calculated with a few non-deep-learning models introduced in section 2.1.1. These are Lagrangian Persistence (extrapolation nowcast) and LINDA-D, the deterministic version of LINDA.

Concerning baseline models other than RainNet, before all predictions, the bounding box is applied, input data is converted to rainrate using Eq. 3.1 and thresholded at 0.1 mm/h and possible NaN values are set to zero. After predictions are made, these are thresholded at 0.1 mm/h and converted to back to reflectivities by first using Eq. 2.1 with parameters  $A$  and  $b$  stated in section 3.2 and then converting Z to dBZ for storage. When predictions are read for metric calculation, they are once again converted to rainrate before metric calculation.

Both extrapolation and LINDA-D, as well as models described in section 3.4.2 use Lucas-Kanade optical flow for advection field estimation with default parameters and four input frames with field interpolation. These models also perform advection using Pysteps' Semi-Lagrangian backward scheme, using cubic interpolation, again with other parameters left to default values. LINDA-D has stochastic perturbations set off with the `add_perturbations` flag and otherwise uses default parameters.

### 3.4.2 Baseline Probabilistic Models

Probabilistic skill verification was performed against STEPS and LINDA-P. The former was chosen as it is a well-established and broadly used model offering good performance for a relatively low inference cost of 1-2 minutes to produce a 3 hours nowcast of 24-48 ensemble members on a modern CPU for a domain of 512x512 pixels [45]. It thus makes sense to want BCNN to perform at least on the same level as STEPS, considering that the time required for inference with it might not be much lower than that. LINDA-P, also known as the probabilistic variant of LINDA, is computationally more

expensive but represents the state-of-the-art in terms of probabilistically predicting localized high-intensity rainfall. Consequently, LINDA-P makes for an interesting benchmark for those cases.

The same preprocessing and postprocessing pipelines as for deterministic cases (section 3.4.1) are used with baseline probabilistic models. Also, the same optical flow and extrapolation parameters are used, and again, Pysteps 1.6.1 default values for methods are used unless stated otherwise.

As stated in section 3.1, probabilistic nowcasts produce 48 ensemble members. STEPS has parameters set to match the input data, that is `km_per_pixel` is 2.0 and `timestep` is 5.0, also `R_thr` is set to -10. LINDA-P has the same parameters as LINDA-D, except that naturally stochastic perturbations are set to True with the `add_perturbations` flag.

## 3.5 Verification Visualizations and Metrics

### 3.5.1 Deterministic Prediction skill evaluation metrics

Deterministic prediction skill scores were calculated in order to compare the raw predictive skill of ensemble nowcasts to the skill of deterministic models. Such comparison is an important facet of verification as low skill in deterministic scores would even for an otherwise competent model mean that it would benefit from being complemented by a stronger deterministic model. In the case of ensemble nowcasting, Deterministic scores were calculated for ensemble means. Implemented deterministic metrics are divided into four different categories, roughly complementing each others. These are continuous, categorical, and spatial scores, as well as radially-averaged power spectral density.

Continuous scores scores are as their name suggests distance metrics used for the evaluation of continuously valued predictions, i.e. regression tasks. The continuous scores used are Mean Error (ME) and Mean Absolute Error (MAE). ME is defined as

$$\text{ME} = \frac{1}{N} \sum_{i=1}^N y_i - \hat{y}_i, \quad (3.4)$$

where we sum over the  $i = 1, \dots, N$  pixels in the radar image,  $y$  denotes ground truth and  $\hat{y}$  the prediction made. The principal utility of Mean Error is detection whether predictions are biased towards too low or too high rainrates at a certain point in time. On the other hand,

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|, \quad (3.5)$$

only cares about the magnitude of the errors in predictions by taking the absolute value, so it gives an idea of the prediction skill over the images *on average*. Nevertheless, this is not enough to accurately assess the skill of a nowcast method, because not all pixels are equally important, as most of them have no rain or very low rainrates that are not interesting to predict. Additionally, A poor forecast may have good MAE and vice versa. To illustrate this, a forecast failing to predict localized intense rainfall, but otherwise accurately capturing light rainfall over large areas will usually have a small MAE but will have low operational usefulness.

This limited utility of continuous scores serves as a motivation to introduce categorical scores. These are based on the principle of comparing the presence or absence of a rain event in observations and predictions as defined by having a pixel exceeding a threshold value  $R_{\text{THR}}$ . The categorical scores are defined by dividing events into four categories, that are true positives (TP), i.e rain events that were correctly predicted, true negatives (TN), i.e. lack of rain event that was correctly predicted, false negatives (FN), i.e. rain that wasn't successfully predicted, and false positives (FP), i.e rain that was erroneously predicted. These categories and their relations are illustrated in Table 3.1. Scores derived from those quantities that were used are probability of detection (POD) [53] defined as

$$\text{POD} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (3.6)$$

, which simply tells the probability that an event really occurring is correctly predicted. Next, false alarm rate (FAR) [53] is calculated as

$$\text{FAR} = \frac{\text{FP}}{\text{TP} + \text{FP}} \quad (3.7)$$

<i>Precipitation is ...</i>	<b>Observed</b>	<b>Not observed</b>
<b>Predicted</b>	TP (Hits)	FP (False Alarms)
<b>Not predicted</b>	FN (Misses)	TN

Table 3.1: Contingency table of categories used. Alternative nomenclature is shown in parentheses.

which reversely indicates the percentage of positively predicted events not actually happening. Critical success index (CSI) [53], which is defined as

$$\text{CSI} = \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}} \quad (3.8)$$

is computed and aims to generally assess the performance of the forecast by taking the proportion of correct positive event predictions out of critically important cases, that is those excluding true negatives but including both false alarms (FP) and misses (FN). Lastly, the equitable threat score (ETS) [28] is defined as

$$\text{ETS} = \frac{\text{TP} - \text{rnd}}{\text{TP} + \text{FN} + \text{FP} - \text{rnd}}, \quad (3.9)$$

where  $\text{rnd} = \frac{(\text{TP} + \text{FN})(\text{TP} + \text{FP})}{\text{TP} + \text{FN} + \text{FP} + \text{TN}}$

was computed. ETS aims to improve CSI assessment of forecast skill, by attempting to estimate the amount of random TP among the prediction using the term  $\text{rnd}$ , and remove that number of data points from the calculations.

The  $R_{\text{thr}}$  thresholds chosen for the performing verification are 0.5, 1.0, 5.0, 10.0, 20.0, and 30.0 mm/h. 0.5 mm/h corresponds to very light rainfall and should be easy to predict, while higher thresholds like 20.0 and 30.0 mm/h correspond to very heavy rainfall, which are very difficult to predict even for short leadtimes.

In addition to evaluating prediction skill above rainrate thresholds, it is also important to be able to evaluate the nowcast at multiple scales. The reason for this is that bigger scales have more predictability, and so predicting smaller scales is more difficult while also being of high importance in the context of heavy localized rainfall.

As such, spatial verification scores, namely Fraction Skill Score (FSS) [49] and intensity-scale verification using FSS are added to the panoply of metrics used. FSS is a verification metric aiming to estimate the prediction skill above a certain threshold at different spatial scales. It works by calculating a binary threshold exceedance map for forecasts and observations, averaging it over gaussian windows of different lengths representing scales, and calculating for each of those a Mean Squared Error (MSE) skill score relative to a reference low-skill forecast [49]. Calculating spatial verification scores for a matrix of rainrate intensity threshold and spatial scale is known as intensity-scale verification. Such verification was performed using FSS for thresholds of 0.5, 1.0, 5.0, 10.0 and 20.0 and 30.0 mm/h, as well as spatial scales of 1, 2, 4, 8, and 16 km.

Last but not least, the ability to maintain small-scale details as the forecast leadtime increases is related to the skill of the nowcast with regards to the spatial scale. Failure in maintaining those details will result in low skill at small scales and a blurred forecast demonstrated by a dip in nowcast small-frequency power spectral density. Hence the last deterministic verification metric chosen is Radially-Averaged Power Spectral Density (RAPSD). This metric as its name suggests provides a way of calculating power spectral densities for 2D images such as nowcasts, independently of direction. RAPSD characterizes the amount and type of blurring occurring in deep-learning based models verified. It is calculated for 15, 30, 60, and 180 minute leadtimes.

**!FIXME fix leadtimes, threshs, scales, before submit FIXME!**

### 3.5.2 Visual Evaluation of nowcast predictive uncertainty and Skill

In order to visually assess ensemble nowcasts, ensemble mean and standard deviations of nowcasts were plotted and compared to radar observations. In addition, rainrate exceedance probability estimations for the ensemble were plotted for thresholds of 0.5 and 5.0 mm/h. Leadtimes chosen for these visualizations included 5, 15, 30, and 60 minutes or alternatively 5, 60, 120, and 180 minutes for studying longer less skillful prediction time-scales.

These visualizations were made for two distinct cases, one containing purely convective rainfall and the other containing a mix of stratiform and convective rainfall. This heterogeneous set serves to assess differences in the developed method with regards to the rainfall type, knowing that convective events are notoriously harder to predict.

### 3.5.3 Probabilistic Prediction skill evaluation metrics

Deterministic verification scores are not enough to accurately assess ensemble forecast skill, because the advantage brought by ensembles does not lie in their mean value, but rather in the breadth and quality of their distributions, allowing to weight in multiple possible scenarios. Consequently, specialized metrics designed to assess probabilistic forecasts are needed. For this work, four different probabilistic metrics are used for verification. These are the Continuous Ranked Probability Score (CRPS), rank histograms, reliability diagrams, and Receiver operating characteristic (ROC) curves.

CRPS is a metric generalizing MAE for probabilistic forecast. It is defined as

$$\text{CRPS}(F, x) = \int_{-\infty}^{\infty} (F(y) - \mathbb{1}(x \geq y))^2 dy \quad (3.10)$$

where  $F$  is the forecast cumulative density function (CDF) and  $\mathbb{1}(x \geq y)$  is the empirical CDF of the observation  $x$ . CRPS aims thus to represent the distance between those cumulative distributions as a proxy to estimate ensemble forecast skill.

A rank histogram is a measure that from the rank of the radar observation among ensemble members, builds a histogram. For each pixel, the rank is determined and a bin is incremented accordingly. The shape of the histogram is indicative of whether the spread of the ensemble is representative of the true spread of observations. Variability in the ensemble equal to the uncertainty of observations would make a flat histogram. A convex histogram would mean that ensemble spread under-estimates true uncertainty, whereas a concave histogram would mean that uncertainty is over-estimated, that is that the ensemble is more uncertain than it should be. Furthermore, skewness of the histogram gives a hint on whether there is any kind of bias in the predictions of the ensemble. This goes a step further than ME, as it does not simply tell the average error, but estimates its distribution in a sense.

ROC curves [36] are a verification method assessing the ability of a forecast to discern between positive and negative events, that is in the present case pixels with or without exceedance of a particular rainrate threshold. This is accomplished by keeping track of false alarm rates (FAR) against probability of detection (POD) of an event. Probabilities are divided into a certain number of bins over which corresponding FAR are averaged, and a curve is formed. A random forecast corresponding to no skill corresponds to a line, and an increasing area under the curve indicates better discernment ability.

A reliability diagram presents observed relative frequencies of events (rain-rates exceeding a certain threshold value) against their average forecast probabilities. When these two values are close to each other, the forecast is said to be reliable. This means that a reliable forecast is defined as having its reliability diagram close to the straight line corresponding to observed relative frequencies and forecast probabilities output by the model being equal. In practice, reliability diagrams are built by dividing forecast probabilities into bins with associated observed relative frequencies. The diagram is often accompanied by a histogram of event counts in those bins called the sharpness histogram. A model's forecasts can be said to be sharp if most predictions are close to zero or one probabilities. Sharpness can be seen as a measure of the "decisiveness" of a forecast.

ROC curves are conditioned on observation of an event, while reliability diagrams are conditioned on its forecast. Because of this they complement each others well in the evaluation of ensemble prediction skill. For probabilistic metrics, the same rain intensity thresholds as in deterministic metrics, namely 0.5, 1.0, 5.0, 10.0, 20.0, and 30.0 mm/h are used. CRPS is calculated for the whole 36 leadtimes covering 3 hours, while other metrics are calculated for leadtimes of 15, 30, 60, 120, and 180 minutes.

**FIXME Fix LTs, thresholds at the end, after experiments are finished** **FIXME!**

All of the scores described in sections 3.5.1 and 3.5.3 complement each others, and forecast skill is thus estimated as a combination of them, as no score is able to capture all of the facets of a great nowcast.

# Chapter 4

## Results

This chapter will comprehensively present the results obtained through the experiments performed. First, the results of the model hyperparameter optimization process are shown. Then, a nowcasting example is given for two meteorologically interesting cases for one selected model, with other versions and an existing baseline given in the appendix. Lastly, both deterministic and probabilistic performance of the model variants against baseline is presented.

### 4.1 BCNN Hyperparameter Optimization

Gaussian NLL loss data uncertainty parameter  $\sigma$  was chosen to be  $1e^{-4}$  amongst choices of  $1e^{-2}$ ,  $1e^{-3}$ , and  $1e^{-4}$ . As gaussian NLL loss varies in magnitude when varying  $\sigma$ , the choice was done using by observing the ETS skill score at different thresholds on leadtimes from 5 to 30 minutes on a deterministic RainNet, with the scores depicted in Figure A.1. Smaller  $\sigma$  values appear to lead to sharper forecasts, while bigger ones to blurrier forecasts.

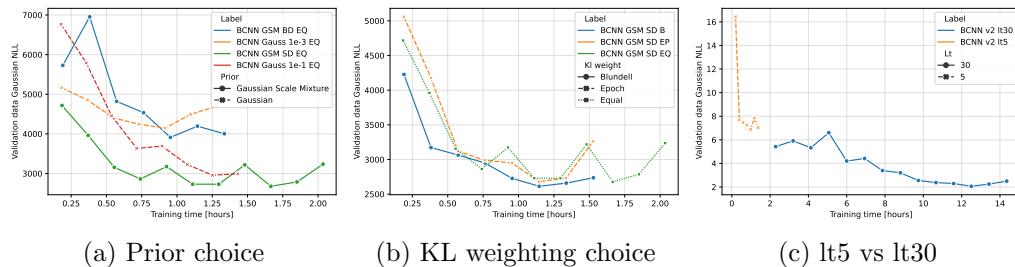


Figure 4.1: The validation set Gaussian NLL loss of different BCNN candidates against total training time, highlighting the model selection process.

The first hyperparameter to be optimized for the Bayesian CNN is the prior. Here, it was chosen based on validation Gaussian NLL loss at model convergence. To give an idea about the natural distribution of network parameters under no regularizing prior, the parameter distribution for the deterministic RainNet with  $\sigma = 1e^{-4}$  is shown in Figure 4.2. distribution is centered at 0 and has high kurtosis, mostly containing parameter values between  $-0.2$  and  $0.2$ .

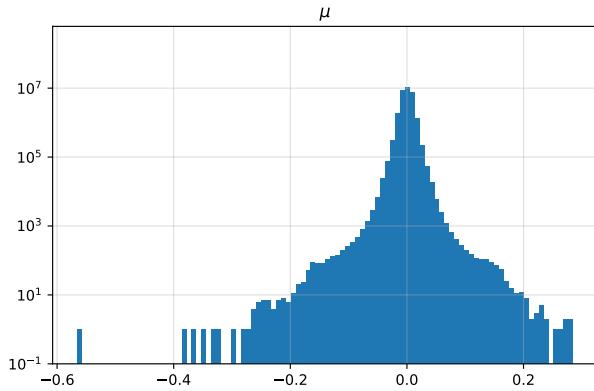


Figure 4.2: RainNet (lt5 with  $\sigma = 1e^{-4}$ ) parameter value distribution after training. Histogram with 100 bins.

Four priors were used as candidates based on knowledge of the "natural" parameter distribution. Two are Gaussian  $\sim \mathcal{N}(0, 1e^{-1})$  (1) and  $\sim \mathcal{N}(0, 1e^{-3})$  (2); and two are Gaussian Scale Mixtures with  $\sigma_1, \sigma_2 = 1e^{-1}, 1e^{-3}$  (3) and  $\sigma_1, \sigma_2 = 3e^{-1}, 3e^{-2}$  (4). Prior (1) serves as a weak prior, allowing most "natural" parameter values. Prior (2) on the other hand is much stricter, restricting the range of possible values to be much smaller than it would be without it. Prior (3) mixes (1) and (2) to try to figure out if the approach of scale mixture has benefits. Prior (4) finally is based on the observation that the empirical distribution of deterministic RainNets is relatively close the PDF of this prior, making it interesting to look at whether using it results in performance improvements.

The results of the optimization are shown in Figure 4.1 (a). Here the KL-weighting scheme is set to equal weighting and the likelihood  $\sigma = 1e^{-4}$ . It is shown that best performance is achieved by prior (4), making it the choice to be used in later experiments. The weaker prior (1) is slower to converge but achieves relatively good asymptotic performance, while the stricter prior (2) gives decent results more quickly but has trouble converging to low loss values. This was substantiated by prior (2) making blurrier

forecasts. Prior (3) behaved as a middle-ground between priors (1) and (2). It benefited from the same fast initial convergence as prior (2) but achieved asymptotic performance somewhere between the two.

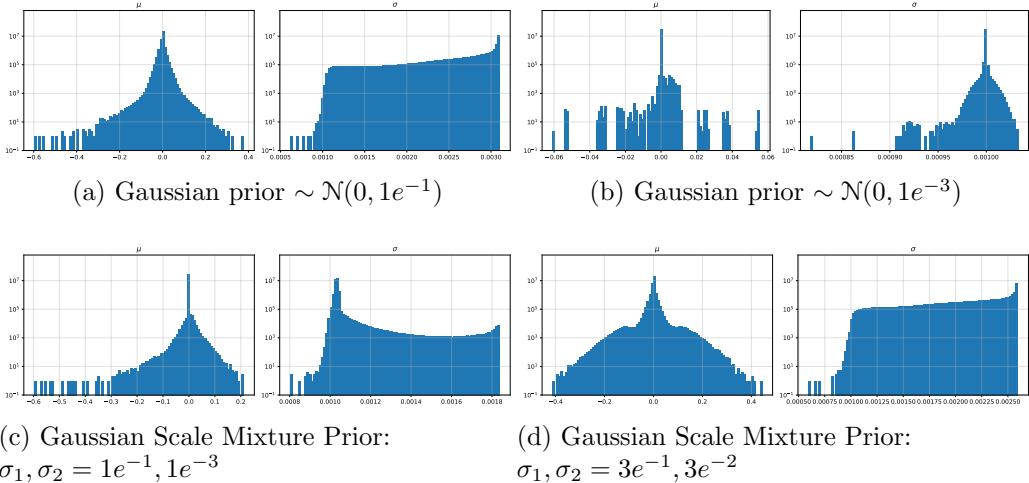


Figure 4.3: The relation of the initial prior distribution chosen to the final parameter posterior distributions across the network.  $\mu$  histograms describe the parameter mean distributions, whereas  $\sigma$  distributions describe the distribution of their standard deviations. The KL-weighting scheme here is set to equal weighting.

The effect of the prior on the parameter mean and STD distributions is portrayed in Figure 4.3. It is seen that for all priors except (2), STD distributions end up tilted towards bigger values than the initial scale  $1e^{-3}$ , pointing to a possibility that it was way too small and that combined with a small learning rate ( $1e^{-4}$ ), there wasn't just enough time for them to grow where they should have. Concerning means, it is observed that the breadth of the posterior mean distribution is closely associated to the prior distribution breadth. With prior (2), it is seen that, most means are actually very close to 0 and while scales are close to their initial values, the biggest deviations tend to lower them, which can be intuitively understood as the prior scale is the same as the initial posterior scale.

With the prior chosen, the next hyperparameter to be optimized is the KL weighting scheme. The results are depicted in Figure 4.1(b). Best performance is observed with the Blundell weighting scheme which was chosen for further experiments, with close performance achieved using equal weighting. Although relatively close in performance, the Epoch weighting scheme did not deliver in the end, although it showed promising results in preliminary

experiments, where the Blundell scheme did not work. The effect of the KL weighting scheme on the weight distribution is shown in Figure A.2. Here, the Blundell scheme differs considerably compared to others in terms of scale distribution, as they are more closely centered around the initial value and the shift is directed towards smaller values.

Finally training times for lt5 and lt30 variants are compared in Figure 4.1(c). The training is shown for the model BCNN lt5 V2, which is an additional experiment not used for other verification. It is seen that continuing the training with 30 minute leadtime (lt30: 6 predictions) requires much more computational resources than simply using 5 minute leadtime (lt5: 1 prediction), which is only worth if skill is improved as a result. In the current example, using a likelihood  $\sigma = 2e^{-2}$ , likelihood loss is improved with lt30, but the opposite is found using likelihood  $\sigma = 1e^{-4}$ .

## 4.2 Nowcasting Examples

Here, example nowcasts made with the model developed are shown for two distinct events. The main three models retained for verification are The main model showcased is BCNN 1t5, with additional models showcases in Appendix A.2. Predictive mean and uncertainty, as well as exceedance probabilities for 0.5, 5.0, and 20.0 mm/h are shown.

### 4.2.1 Case study 1 : Rapidly Moving Mixed Stratiform and Convective Rain Event

The predictive mean and uncertainty for BCNN 1t5 on Case 1 is shown in Figure 4.4. It is observed that mean predictions get expectedly increasingly smoothed out with time, and that predictive uncertainty is highly correlated with predicted mean rainrate. Additionally, predictive uncertainty increases over time, achieving top values exceeding 50 mm/h in areas of predicted heavy rain at one-hour leadtimes. It is seen that advection and the overall movement of the front is well-modeled, while smaller-scale details even if corresponding to medium or heavy rain, are quickly lost.

Smoothing out of predictions has an important effect of reducing model capacity to predict high rainrates, which are often concentrated over small areas. This is an especially important effect because smoothing does not originate simply from averaging, but is also present in individual samples, due to the combined effect of the loss function and of the iterative prediction scheme, as described in section 3.3.1.

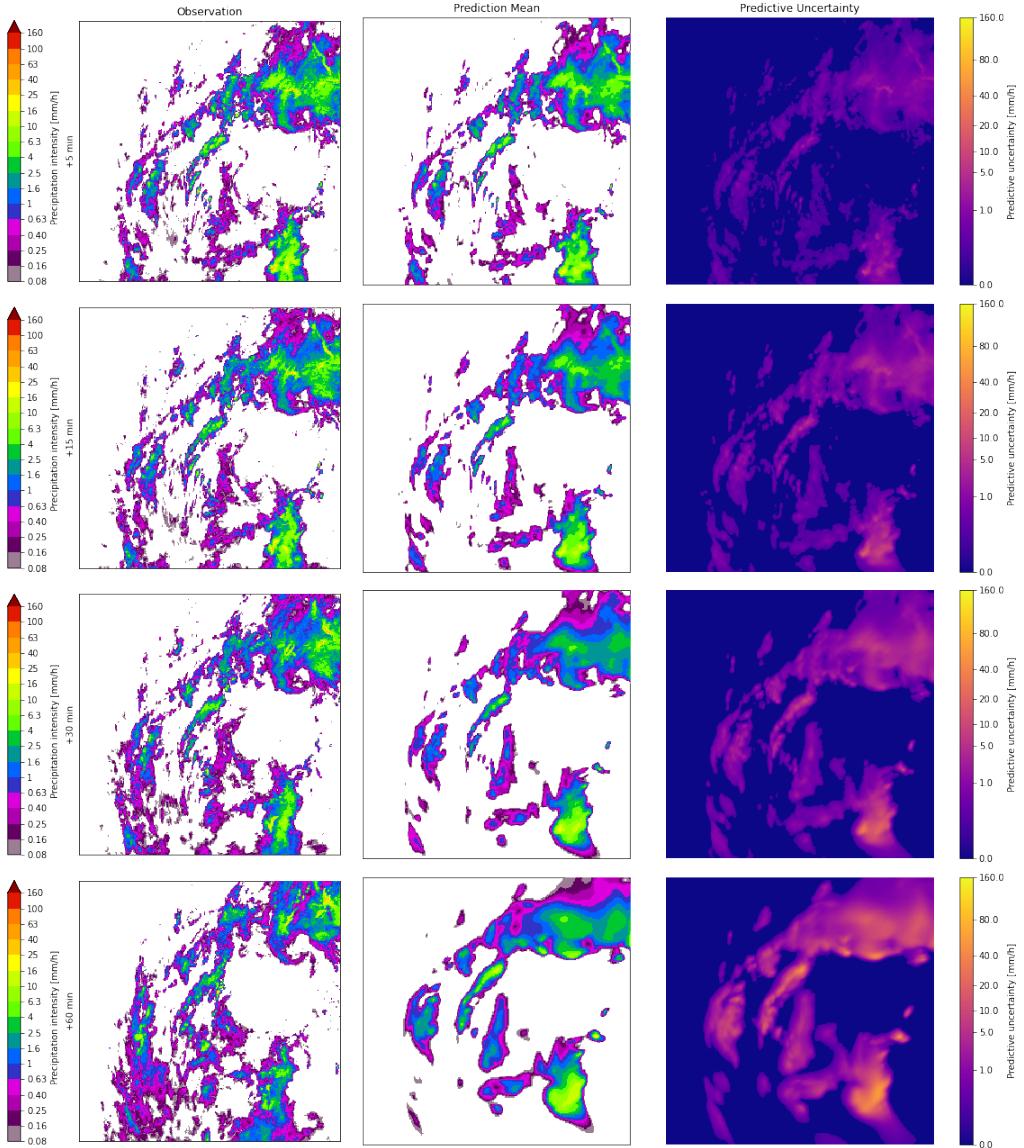


Figure 4.4: Predictive mean and uncertainty (two standard deviations) for Case 1 on the 25th of May 2019 starting at 1PM for the BCNN 1t5 model.

Exceedance probabilities for BCNN 1t5 on Case 1 are shown in Figure 4.5. Both for this model and in general, it is seen that exceedance probabilities are lower for high leadtimes and high thresholds. In the case of this model, all exceedance probabilities seem to be smoothed over time, which is indicative of the uncertainty of the extent of rain being at least to some degree taken into account. It is also seen that especially at 20.0 mm/h, the predicted exceedance probability at one-hour leadtime does not always originate from

current time high rainrates. This is seen from the top of the image at +5min, where a small exceedance probability is present, but disappears at further leadtimes, before reappearing independently.

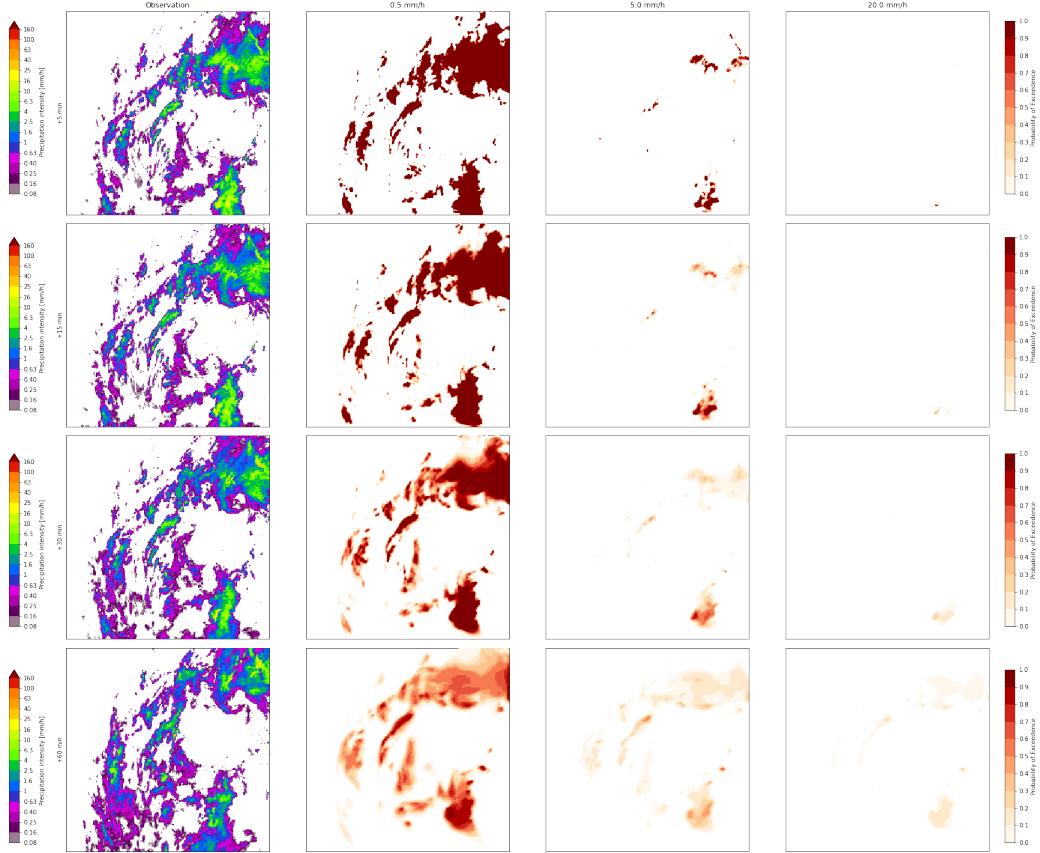


Figure 4.5: 0.5, 5.0, and 20.0 mm/h exceedance probabilities for Case 1 on the 25th of May 2019 starting at 1PM for the BCNN 1t5 model.

Both in the plot of prediction mean and uncertainty, as well as in that of exceedance probabilities, it is seen that because the rain is advected towards the north, lower parts of the image can not be predicted.

#### 4.2.2 Case study 2 : Mixed Stratiform and Convective Rain Event with Spinning Motion

The predictive mean and uncertainty for BCNN 1t5 on Case 2 is shown in Figure 4.6. In this case, the same overall features are observed as in Case 1. Here, the spinning motion is correctly predicted and the mean rainrates

are relatively close to ground truth. As in Case 1 but especially here, many small features are lost and especially small convective cell evolution is not correctly predicted.

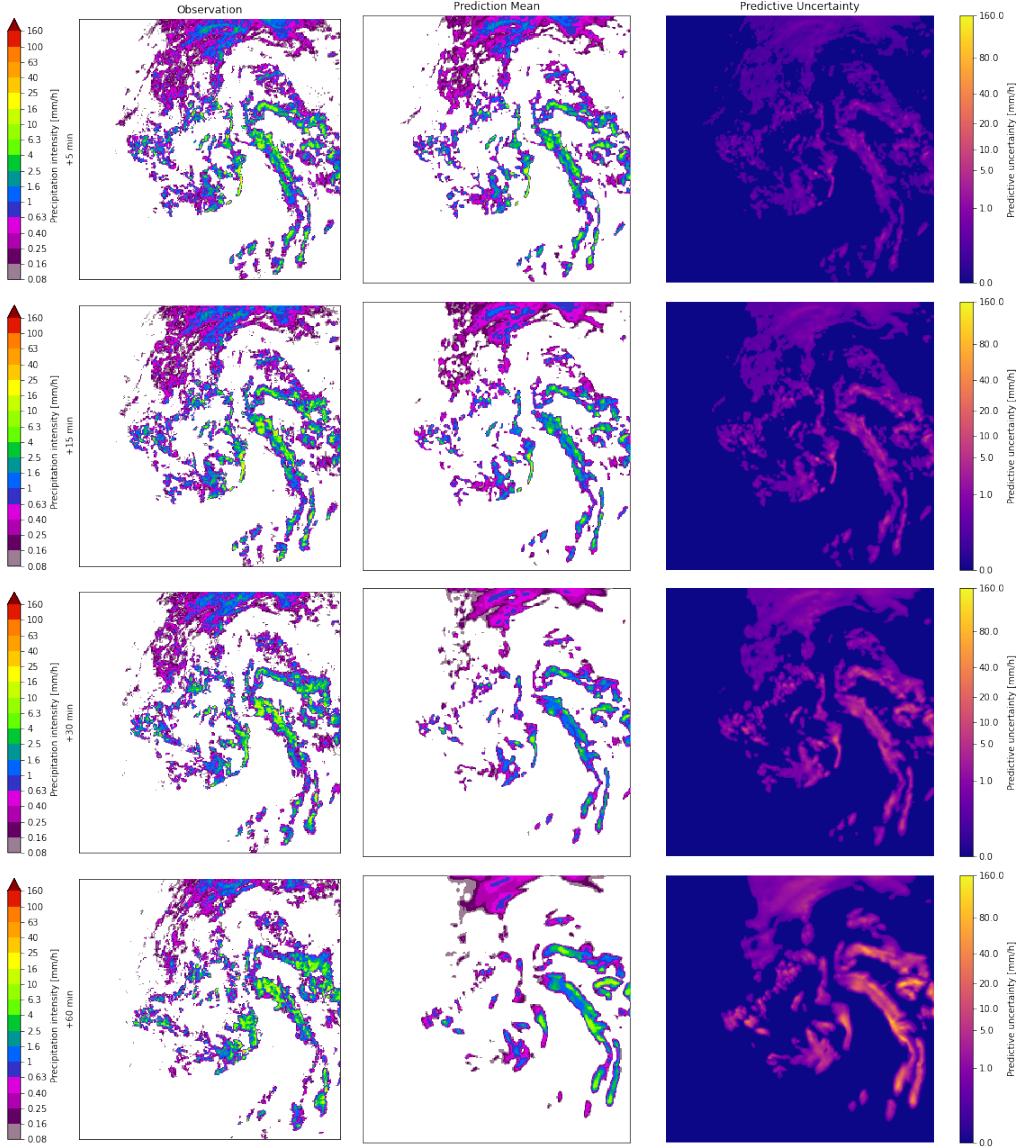


Figure 4.6: Predictive mean and uncertainty (two standard deviations) case 2 2021-05-18 21:00:00

The exceedance probabilities are shown in Figure 4.7. Many things happening here are similar as in Case 1, but especially striking is lower exceedance probabilities at one-hour leadtime for the 0.5 mm/h threshold.

Here, in many samples the spinning "arms" were lost by that time because of their thinness and the tendency of thin or small features to disappear.

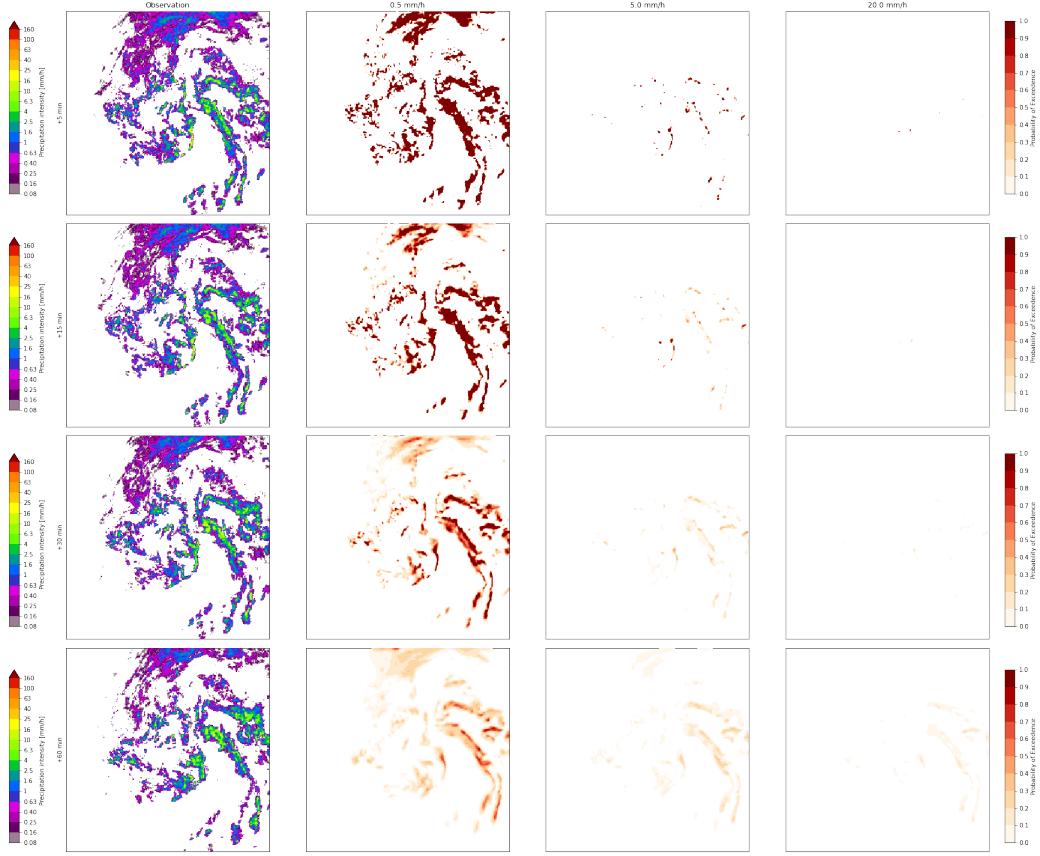


Figure 4.7: 0.5, 5.0, and 20.0 mm/h exceedance probabilities for Case 2 on the 18th of May 2021 starting at 9PM for the BCNN 1t5 model.

### 4.2.3 Additional Case Study 1 results

Case 1 prediction mean and uncertainty, as well as exceedance probabilities are respectively plotted for BCNN 1t5 new in Figures A.3 and A.4, BCNN 1t30 in Figures A.5 and A.6, and finally STEPS in Figures A.7 and A.8. For BCNN 1t5 new, results are similar as in the BCNN 1t5 case, but exhibit a much bigger bias towards predicting strong rain, with more ensemble members predicting increasingly high rainrates. From exceedance probabilities, this can be seen as overall higher values, making the model more susceptible to predicting stronger rainrates.

For BCNN 1t30, it can be seen from both figures that the overall skill is lower than for 1t5 cases. In particular, predictions are more smoothed out and more details are lost, hinting to the conclusion that 1t30 training did not work as expected. Here, predictive uncertainty does not increase with leadtime, and higher rainrates are often not predicted at all. However it can be seen that the network predicts consistently high exceedance probabilities, and that those are not (esp. 0.5 mm/h) smoothed out spatially with increased leadtime.

Finally, looking out at Figures of the STEPS baseline, it is easily seen that despite being a simpler model, it predicts better spatial uncertainty in future rain. However STEPS has a hard time predicting exceedance probabilities for high (esp. 20.0 mm/h) thresholds. One striking feature is also that in the case of STEPS, predictive uncertainty is more spread out than predictive mean, which is a behaviour not observed in any model developed here.

### 4.3 Deterministic prediction skill (Metrics)

The following results describe the skill of ensemble means of developed models against baseline models. LINDA-P ensemble means are omitted for clarity because their skill was almost identical to LINDA-D.

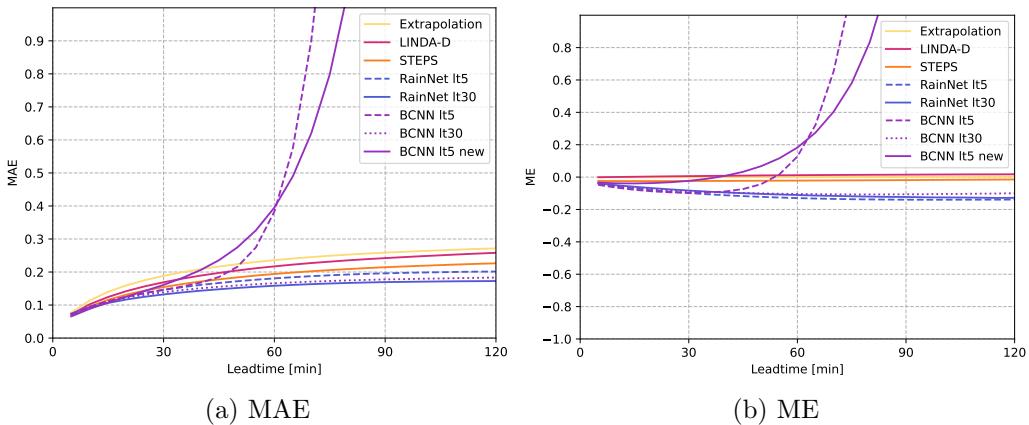


Figure 4.8: Continuous Metrics (MAE and ME) for BCNN variants and baselines, up until two-hour leadtime.

Starting off with the continuous metrics, their values are shown in Figure 4.8. BCNN and RainNet models have similar MAE and ME, except for BCNN 1t5 and BCNN 1t5 new cases, where ensemble means exhibited clear bias towards too strong rainrate predictions after 30 minutes. These models shall

be called "positively biased". ME was around zero for classical models and tended lower with other RainNet and BCNN models, showing bias towards too weak predictions ("negatively biased"), mainly due to smoothing. Still, BCNN and RainNet gave mostly clearly superior MAE compared to classical models, which is understandable as Gaussian NLL is a loss function part of the same class as MAE, which tends to optimize for minimizing pixel-per-pixel deviations and by extension MAE.

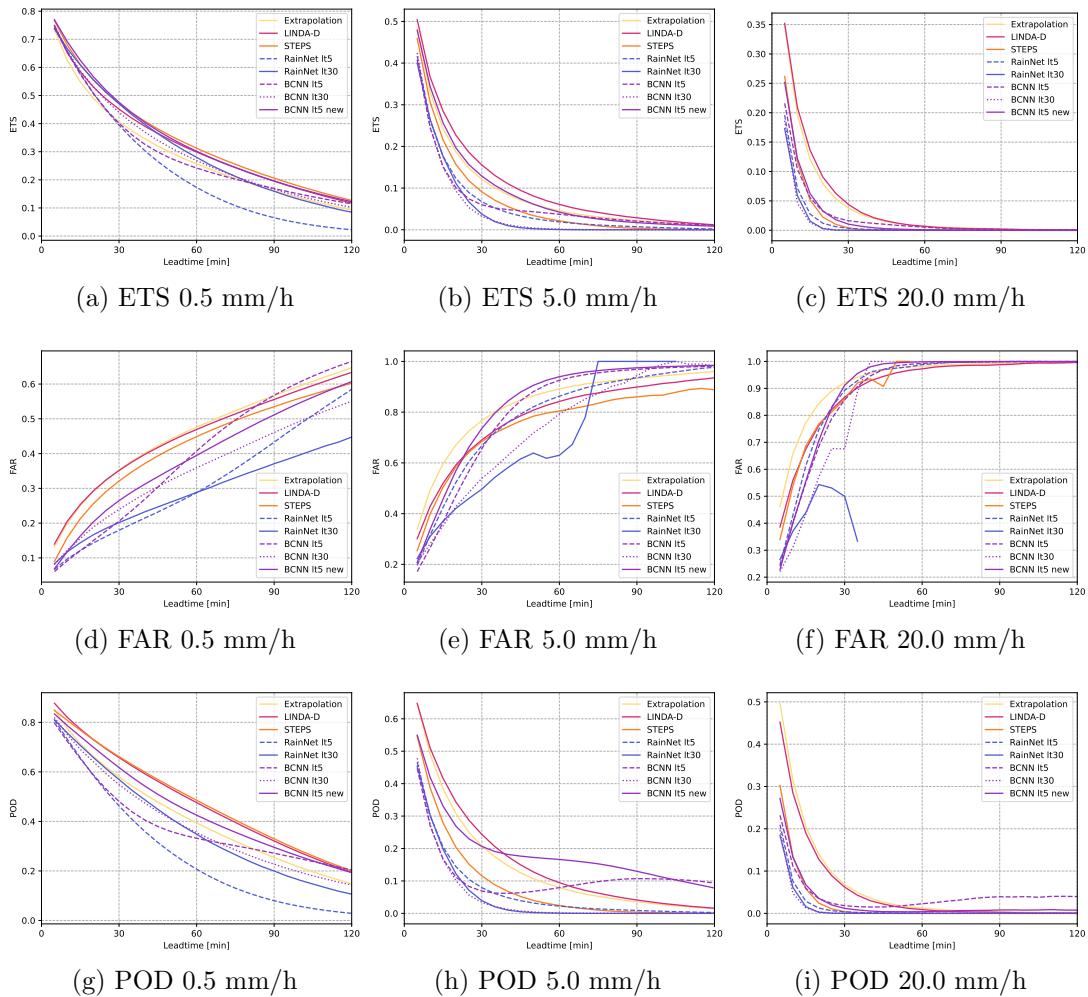


Figure 4.9: Categorical metrics (ETS, FAR, and POD) for BCNN variants and baselines, up until two-hour leadtime at thresholds of 0.5, 5.0, and 20.0 mm/h.

POD, FAR, and ETS metric values are shown in Figure ???. For these metrics, the strongest overall result among developed models is achieved by

**BCNN 1t5 new.** This is not unexpected, as positively biased models tend to perform better with categorical metrics than negatively biased. For POD, a resurgence at further leadtimes is seen in BCNN 1t5 and BCNN 1t5 new models, because of the phenomenon of "exploding rainrates". In other cases, POD is better for classical baselines than CNN based models, when negatively biased. Accompanying high POD of "exploding rainrates" of positively biased models is a very high FAR, demonstrating that this high POD is not necessarily indicative of skill. FAR is interestingly worse for extrapolation based baselines than CNN based models at 0.5 mm/h threshold, a phenomenon that disappears at higher thresholds.

ETS summarizes the categorical skill of the models. Especially at 5.0 and 20.0 mm/h, LINDA-D reigns king of predictive skill, while CNN based models and STEPS lag far behind. Overall, no conclusion can be made about whether the bayesian extension improved the categorical predictive skill, as the only skillful outlier BCNN 1t5 new didn't have a RainNet model with a corresponding training routine to compare to. CSI scores are shown in Figure A.9, but they do not differ much from ETS.

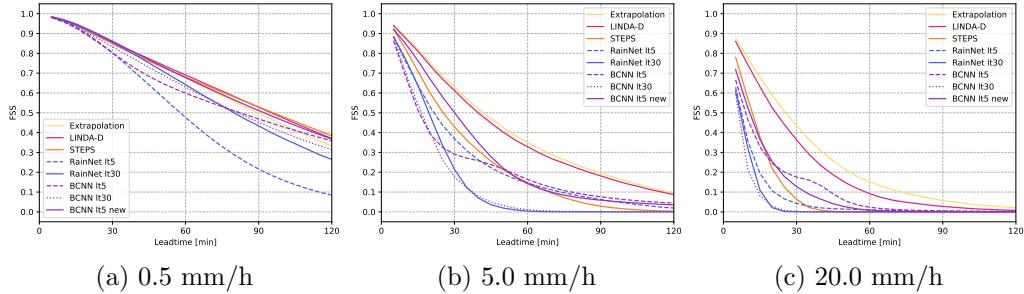


Figure 4.10: FSS metric at a 16km scale for BCNN variants and baselines, up until two-hour leadtime at thresholds of 0.5, 5.0, and 20.0 mm/h.

The FSS metric results at a 16km spatial scale for the thresholds of interest are shown in Figure 4.10. Here, it is found that the difference between Extrapolation based deterministic models and others is strikingly big, with the former being much more skillful at that level of detail. FSS at a scale of 4km is shown in Figure A.10, but the results do not differ much from ETS and CSI at the original scale of two kilometers, as the change is only two-fold.

RAPSD is shown in Figure 4.11 for leadtimes of 15, 60, and 120 minutes. As expected, Extrapolation nowcast which preserves details has a PSD very close to observations, while other models show drop in power at short wavelengths, with steeper drops at further leadtime. Interestingly, this drop in power is bigger at 60 minutes than 120 minutes, which might have something

to do with the OR-masking of NaN values, reducing significantly the area used for calculation at two hours. The biggest drops in small-scale details are experienced by STEPS and BCNN 1t30 as well as RainNet 1t30 models. All 1t5 CNN models experienced less loss in details compared to 1t30 models. It is difficult to conclude whether taking ensemble means had any significant effect on the loss of small-scale details of CNN approaches without looking at individual ensemble members and ditching the OR-masking.

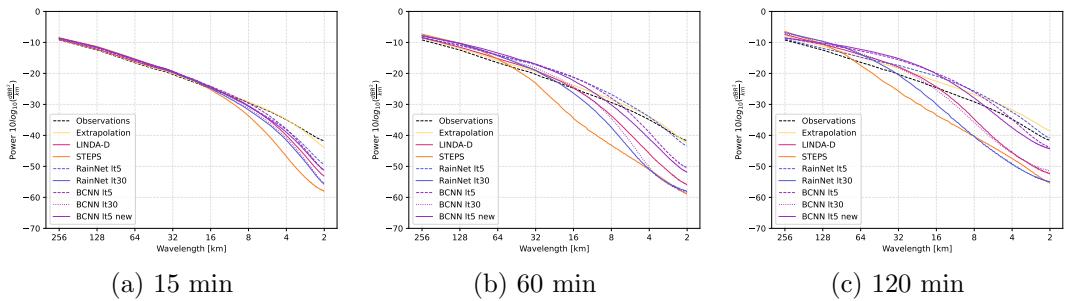


Figure 4.11: Radially-averaged power spectral density (RAPSD) for BCNN variants and baselines, at leadtimes of 15 minutes, one hour, and two hours.

## 4.4 Probabilistic prediction skill (Metrics)

The following section shows probabilistic prediction skill assessment results of BCNN models against STEPS and LINDA-P baselines.

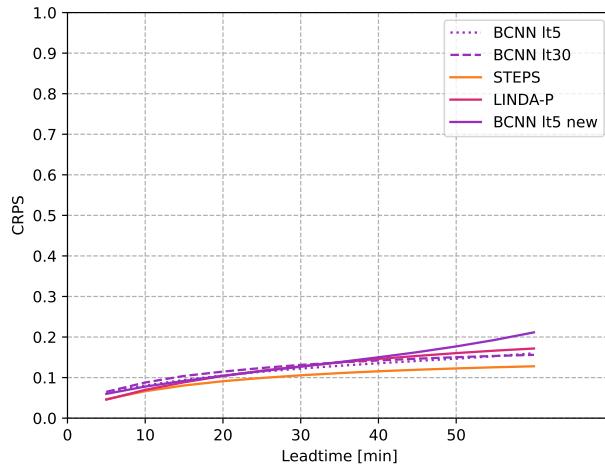


Figure 4.12: CRPS metric for all ensemble models at leadtimes up until one hour.

The continuously ranked probability score for all ensemble models is shown in Figure 4.12 for leadtimes up until one hour. Initially the best CRPS is achieved by LINDA-P, but after 10 minutes STEPS becomes the best. BCNN models consistently perform worse than STEPS, and only LINDA-P gets worse than BCNN 1t5 and BCNN 1t30 after 45 minutes.

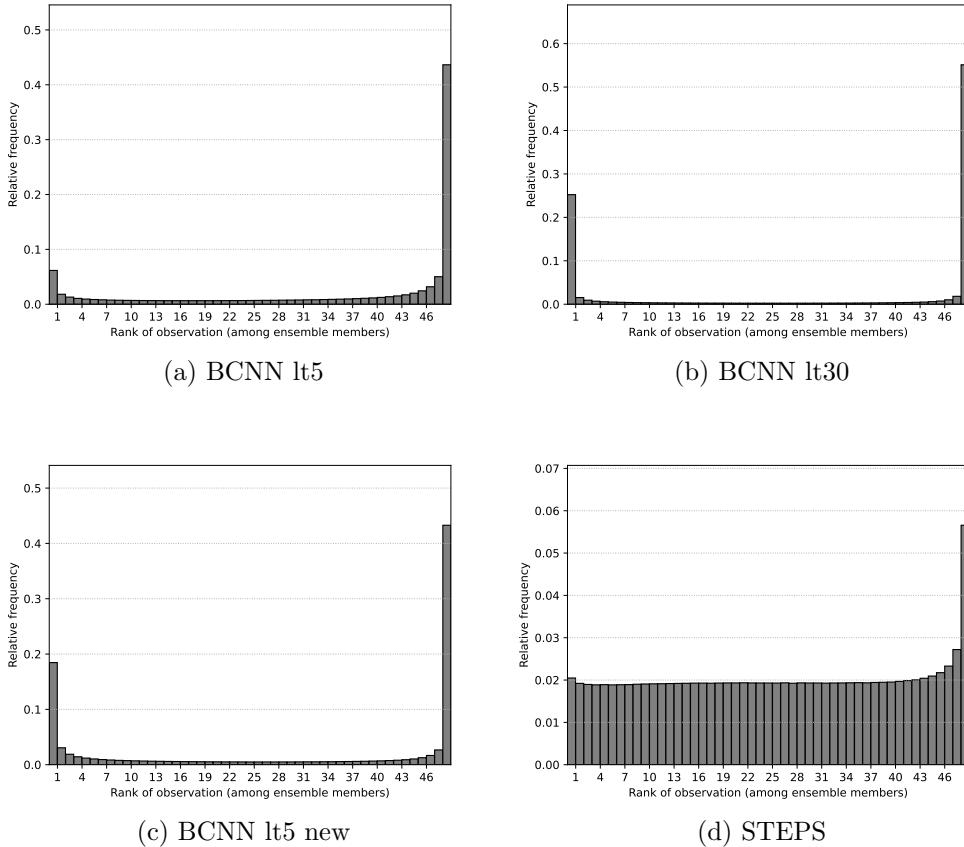


Figure 4.13: Rank histograms for BCNN and STEPS model nowcasts at a one-hour leadtime, for precipitation events exceeding 0.1 mm/h.

Rank histograms for BCNN and STEPS model nowcasts at one hour are shown in Figure 4.13, and the corresponding plot for LINDA-P is depicted in Figure A.11. It is seen that LINDA-P and especially STEPS have very flat histograms, which is indicative of well-calibrated uncertainty, i.e. no under- or over-estimation of true uncertainty is made, as observations are ranked with close to equal probability amongst ensemble members. Only a slight "negative bias" towards weak rain forecasts can be detected by observing

that the rightmost bin is slightly over-represented. That is cases where the observation had the highest rainrate compared to all ensemble members.

Compared to those baselines, BCNN models had all strongly convex histograms, severely under-estimating the true uncertainty of data. Furthermore, BCNN models had the rightmost bin over-represented too, with 40 to 50% of the time, the observation being the highest rainrate. Additionally, BCNN 1t5 new and BCNN 1t30 had around 20% of observations where the observation was smaller than any ensemble member, hinting to many cases where the network is vastly over-confident about its high rainrate prediction.

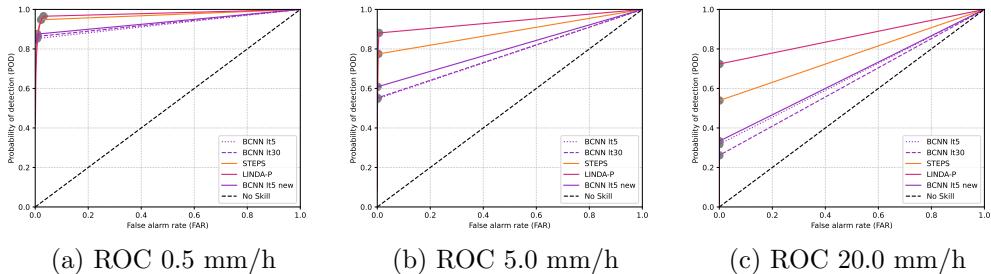


Figure 4.14: ROC curves for one-hour exceedance probabilities at 0.5, 5.0, and 20.0 mm/h precipitation thresholds.

The ROC curve for models at one hour for exceedance probabilities at thresholds of interest is shown in Figure 4.14, and an overview of discrimination power as a function of leadtime and threshold proxied by the ROC area under the curve (AUC) for BCNN models and STEPS is shown in Figure 4.15. From the ROC curves, it is seen that best discriminative power is achieved by LINDA-P, overwhelmingly so at high thresholds. Second was STEPS, followed by the three BCNN models lagging far behind. From the AUC plots, it is visible that for all models, leadtime exhibits a stronger effect on discrimination power for higher thresholds. The AUC plots come to confirm that the discrimination power difference between models derived from ROC curves is consistent varying leadtime and threshold, with perhaps the exception of BCNN 1t5 compared to BCNN 1t5 new performing better on higher, but worse on lower thresholds.

Lastly, the reliability diagram accompanied by its sharpness histogram for each model, at the same precipitation thresholds at a one-hour leadtime is shown in Figure 4.16. Starting with the histogram, the sharpness of nowcasts seems to generally increase with the threshold. At 0.5 mm/h, all models have a relatively similar U-shaped sharpness histogram, indicative of moderate sharpness, with BCNN models perhaps a little sharper. Higher thresholds

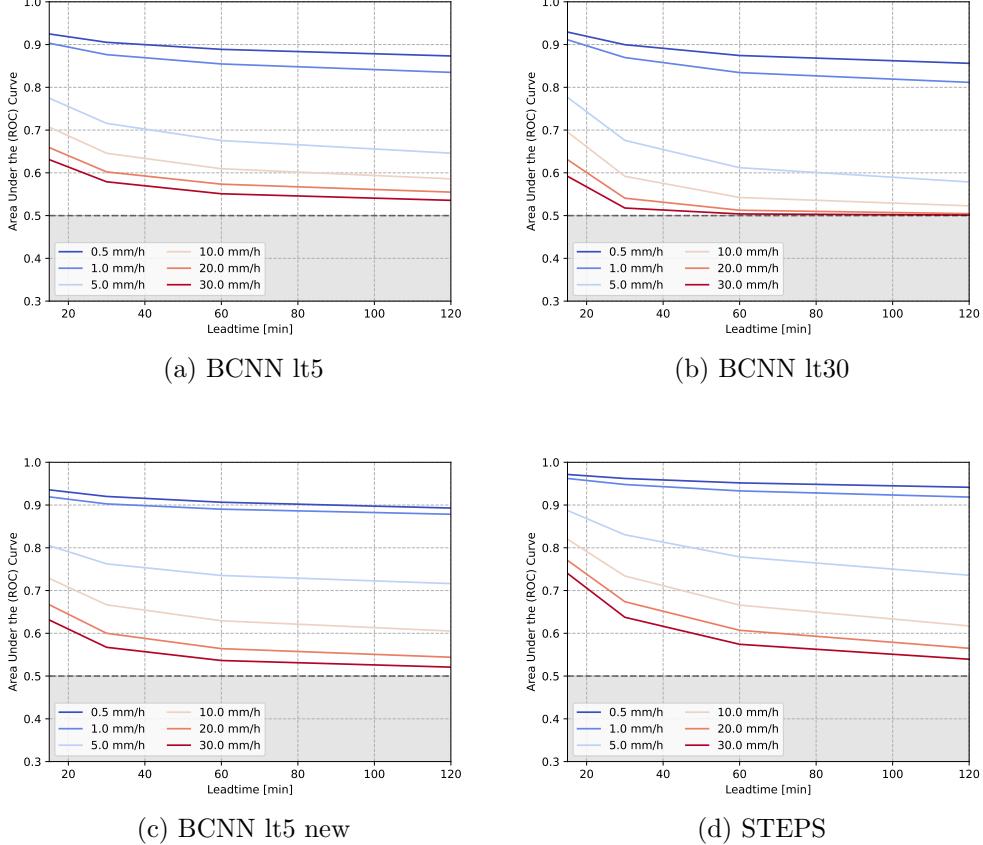


Figure 4.15: ROC Area under the curves summarized for BCNN and STEPS models. The dashed lines correspond to a random forecasts, and grayed out areas to results worse than that of a random (no skill) forecast.

see all models adopt distributions of exceedance probabilities skewed towards smaller values. Generally it seems to be that it is at least one order of magnitude less common to have occurrences of forecast probabilities close to 0.5 for BCNN models compared to LINDA-P and STEPS.

Now turning towards the reliability diagram, it is clear that nowcasts get less reliable for higher thresholds. One will notice that LINDA-P and STEPS both stay both much closer to the black dashed line, indicator of perfect reliability. BCNN models tend here to assign high forecast probabilities to events that will really take place less often, and conversely assign too low forecast probabilities to events that will take place more often in practice. This expresses a lack of reliability. Furthermore, one can see that the reliability diagram of BCNN models gets non monotonous at 20.0 mm/h. This

is due to too small sample sizes in the middle probabilities.

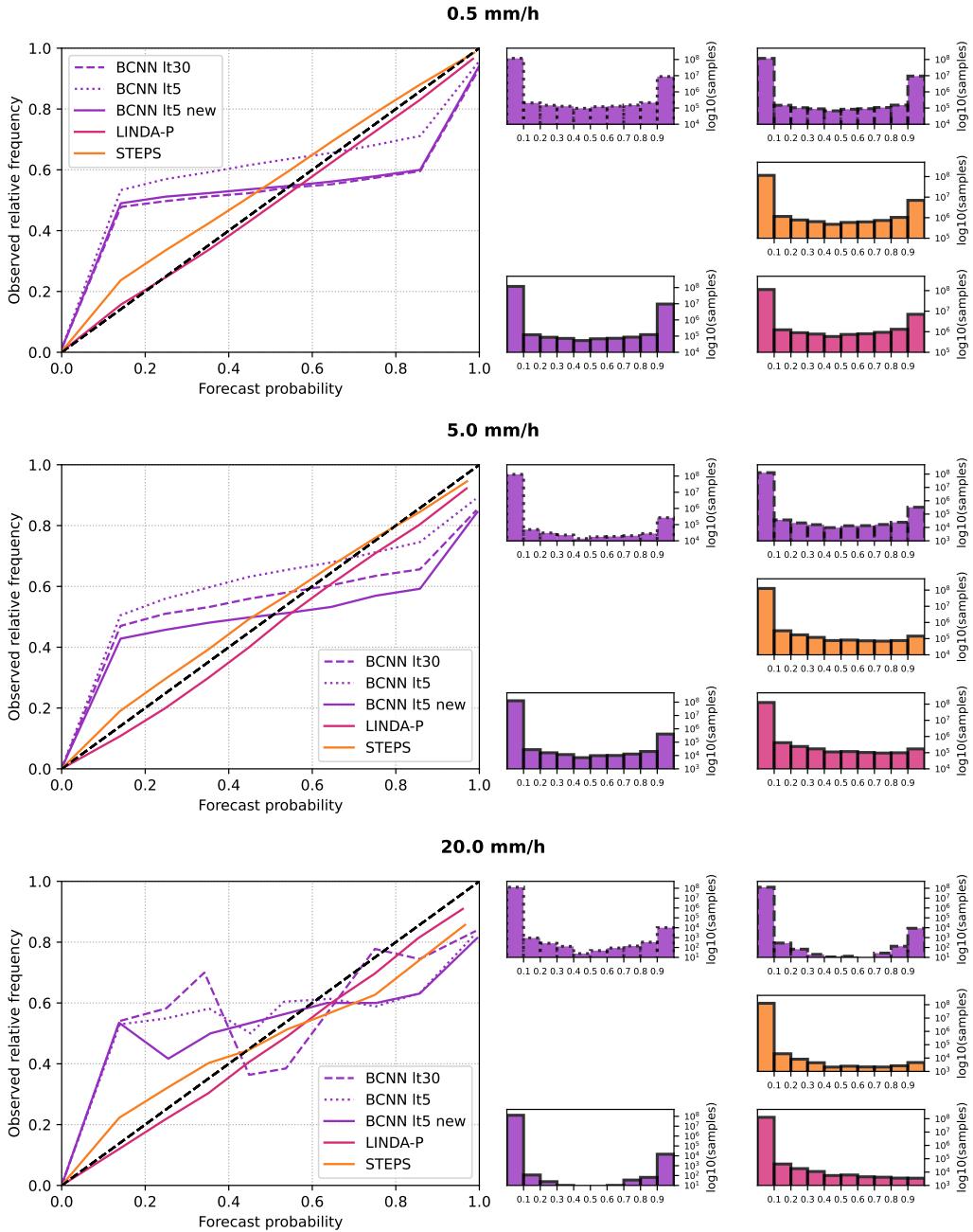


Figure 4.16: Reliability diagrams and sharpness histogram for one-hour exceedance probabilities at 0.5, 5.0, and 20.0 mm/h precipitation thresholds.

# Chapter 5

## Discussion

### 5.1 Goodness of results

BCNN ensemble means produced acceptable deterministic skill, with properties in line with the RainNet baseline [8] and other CNN based models, with characteristic blurring, tendency to be biased or unstable at further leadtimes, difficulties skillfully nowcasting high rainrate thresholds, but high success at lower thresholds and very good MAE. It should be kept in mind that achieving high predictive skill both with low threshold/MAE and high threshold is a difficult problem which is the focus of active research, and is related in the Deep Learning approach to the data imbalance problem of strong precipitation. This means that failure to predict strong precipitation is related to those events representing only a small part of the training data, while being of high importance. Trying to accomodate this by means of event importance weighting is a possibility, but this weighting might easily lead to the network performing worse on all non-extreme cases.

Interestingly, `1t30` training after `1t5` pre-training not only did not improve, but worsened the predictive skill with Gaussian NLL loss, whereas it has in unrelated experiments `1t30` used as a standalone showed much better skill than `1t5` training with LogCosh and MS-SSIM loss functions, owing to a more stable iterative process. One hypothesis why the scheme didn't work out in the present case is that the network should not have been pre-trained with `1t5` but rather directly trained with `1t30`. One related explanation could be that the non-improvement might be due to the constant Gaussian NLL  $\sigma$  data uncertainty parameter. Here with `1t5` training the parameter could indeed be assumed constant for each input, but it can be easily seen why iterative re-plugging in of the outputs breaks that assumption, as more uncertainty is brought into the inputs.

From a probabilistic nowcasting point of view, BCNN did not live up to expectations. Only CRPS is relatively close to baselines, but all other metrics show that BCNN makes disastrously bad predictions compared to STEPS and LINDA-P. Rank histograms show that predictive uncertainty is severely under-estimated, Reliability diagrams that their predicted exceedance probabilities are not very reliable, and thus can not be trusted. ROC curves show that their discriminative power is also considerably worse than baseline, keeping in mind that the ROC curves suffer from too little probability categories, a subject which will be covered below. These results highlight the fact that in its current form, BCNN is not yet a useful probabilistic nowcasting method.

Looking at ensembles and rank histograms it is noticeable that sample variability seems low. This might be related to hyperparameters or more generally to the stochastic model as will be expanded on below, or it might have something to do with the high-frequency detail loss occurring, as quick disappearance of them leads to only slowly evolving large-scale features remaining, which might undermine the variability at further leadtimes. If details were preserved such as in generative model based nowcasting methods, it could help assuming that is indeed part of the problem.

## 5.2 Validity of Assumptions and Experiments

The way that hyperparameter optimization was carried out in this work is a possible source of problems, potentially having led to sub-optimal solutions. This and other validity issues will be covered in this section.

### SVI Issues

When it comes to the Variational Inference method employed, there might have been unguarded assumptions regarding prior and posterior families chosen. In particular, the problem in question might have benefited from a more careful choice of those families based on physical constraints and known behavior of precipitation. Related to this matter, it is known that Variational approximations tend to underestimate uncertainty of learned distributions. [11, 38]. However, it is hard to say if this is the whole story here when observing the rank histograms and reliability diagrams of BCNN models. As the rank histogram shape of them is generally concave compared to models based on stochastic perturbations, the variability of ensembles is indeed less than that of observations most pronouncedly in variational models, as one might expect.

Nevertheless, the magnitude of the effect leaves the possibility that some

parts of the model might have been non-optimal for the task, even within the SVI framework. The reason for thinking this is that in part of preliminary experiments performed, uncertainty estimates for BCNN precursors were way more reliable than in current experiments. *A posteriori*, it was found that one hyperparameter that was different and left un-optimized here is the initialization value of parameter posterior scales, which here was set to  $1e^{-3}$  here and  $1e^{-2}$  in these preliminary experiments. Setting the parameter too high led to difficulties for the model to converge and unstable predicted rainrates, but while it was also known that setting it lower leads to initially smaller uncertainty estimates, no attention was paid while performing experiments to whether this might affect the asymptotic behavior of the estimates. Looking at Figure 4.3, it is seen that scales only deviate from their initial values by a factor of two to three at maximum, hinting the fact that asymptotic marginal uncertainty estimates may very well be connected to the initial scale of parameters.

## Likelihood Cost and Uncertainty

It was also hypothesized whether the magnitude of the Gaussian NLL loss data uncertainty parameter  $\sigma$  influenced the uncertainty estimates, but increasing it from  $1e^{-4}$  to  $2e^{-2}$  as in preliminary experiments had no effect on the estimates. With regards to this data uncertainty parameter, optimizing for it using conventional RainNet was questionable, because the interactions between the complexity and likelihood terms of ELBO are not thus taken into account, and the choice is only based on deterministic skill, which is only a small part and sometimes in contradiction with other aspects of what an ensemble probabilistic model should aim for.  $1e^{-4}$  is also much smaller than the realistic average uncertainty in data in the value range of interest. This uncertainty is expressed in terms of (log-transformed) rainrate (mm/h) and the data has a resolution expressed in reflectivity of 0.5 dBZ, which has for effect that the resolution of the rainrate data is dependent on rain intensity. This calls for non-constant data dependent, e.g. heteroscedastic uncertainties. Another facet of the model making the homoscedasticity assumption questionable is the iterative prediction scheme. This is because the outputs re-fed into the network will inherently have higher aleatoric uncertainty than original inputs, as model (epistemic) uncertainty of previous steps increments data uncertainty of next steps. Also, smaller-scale features often including high-intensity precipitation, have worse predictability and will suffer from higher uncertainty than large-scale features.

## Input Data

Now concerning the input data, the downsampling by a factor of two performed for computational performance reasons might have been detrimental for the validity of the verification of small-scale and especially high-intensity precipitation. There are two factors at play here. One is that because of the scale-dependence of precipitation lifetime, NN-based nowcasting model performance might not scale the same way when changing scale as extrapolation-based classical models, which might make interpreting the performance obtained at bigger scales more difficult. The other factor is that because average pooling is used for downsampling, higher rainrates get diluted, an effect exacerbated by the on average small spatial extent of intense precipitation. Hence, downsampled performance may not necessarily extrapolate to higher spatial resolutions. This is important because in the end, developed precipitation nowcasting models are usually to be used with input data having a grid resolution of 500 meters to one kilometer.

In addition to the downsampling issue, performing the train/validation/test split using six hour blocks might not have been enough to discard all harmful temporal correlations between different time-series, as input images from the end of one block are still correlated to images at the start of the next block, which is part of another split. Thus, even though the correlations are reduced, they do not disappear and the validation and test split performance is still over-estimated for Deep Learning models. A more robust split would have been using non-successive days as units.

## Verification Experiments

RAPSD results suffer from potentially spurious gain in small wavelength power at 120 minutes. This might have something to do with the OR-masking of precipitation reducing too much the area used, or it might be some other unknown effect or mistake in the calculations, which is why RAPSD results are to be taken with a grain of salt.

ROC curves of Figure 4.14 in probabilistic verification suffer from too few probability categories, as the curves exhibit apparent discontinuity, while they should approximately follow a bi-normal distribution. This undermines their accuracy for determining forecast accuracy and resolution.

### 5.3 Directions for Further Work

First and foremost, further work should aim to produce more reliable uncertainty estimates in the current framework. This could be approached by first examining the effect on initial posterior scales on sample variety. Secondly, BCNN would have to be trained and tested using non-downscaled 1km grid resolution reflectivity composites. If these do not resolve most current issues and more time and resources are available, it could be possible to implement other methods for uncertainty estimation using NN and compare them to Bayesian NN with Variational Inference. Different priors, maybe more informed ones or some with varying degrees of strictness, perhaps even non identically distributed among parameters, could be tried to see if that improves predictive skill.

Some other simple tricks that might improve performance of BCNN are pre-learning the parameter means as point estimates and using the Bayesian approach to only learn posterior scales while keeping the means fixed. Assuming a non-biased deterministic model, this could make the learning process much more efficient, raising the quality of uncertainty estimates. Another trick would be to take a pre-trained BCNN model underestimating uncertainty, and add post-processing to the predictions, aiming to calibrate the bias and add stochastic noise for increasing predicted uncertainty. Lastly, despite not improving model performance, weight pruning could be tried for reducing computational expensiveness of making nowcasts. Blundell et al. [12] show that Bayesian Neural Networks with a high number of parameters are very sparse thus and amenable to extensive weight pruning without much loss in skill. Although that study concentrates on fully-connected Neural Networks, it is worth considering if skillful pruning could be efficiently implemented in CNN architectures too, as Shridhar et al. [61] seems to have experienced success doing that.

Despite being one very interesting research direction, decomposition of predictive uncertainty into aleatoric and epistemic uncertainties is not performed, as it was deemed out of scope for this work. Indeed, trying to model the uncertainty the radar image inputs is a separate problem which is not so much of interest in the domain of radar-based precipitation, as so much more can be done to improve performance by ameliorating radar scan fidelity, preprocessing, and choice. Despite of this, one could still view aleatoric uncertainty as a proxy to non-predictability of precipitation given starting conditions represented by the input frames. This could in the case of the iterative models presented at least serve as a metric to quantify the uncertainty created by feeding back in previous model outputs and could

potentially allow the compound aleatoric and epistemic ensemble spread to better match the true distribution of the data. If modeling the uncertainty as data dependent and using heteroscedastic Gaussian likelihood, it would be interesting to take the same approach as Kendall et al. [29] and learn heteroscedastic uncertainties as a separate output channel. This uncertainty learning would eliminate the hairy problem of choosing a homoscedastic data uncertainty value *a priori*, and possibly make for more skillful probabilistic nowcasts. How the separate channel aleatoric uncertainty would be integrated with the ensembles for predictions is still an unsolved problem.

# Chapter 6

## Conclusions

In this work, Bayesian Convolutional Neural Network (BCNN) model variants were developed in attempt to produce skillful and reliable ensemble based probabilistic nowcasts. Two facets of BCNN were evaluated. Firstly the skill of its ensemble means used as a point estimate for forecasting precipitations, and secondly the quality of its uncertainty estimates was evaluated to measure the fitness of the model for probabilistic nowcasting. These facets were compared to existing baseline models, both extrapolation and Neural Network based.

BCNN ensemble averages performed acceptably as deterministic nowcasts, exhibiting many of the known distinctive features of Neural Network based nowcasts, with both good and bad aspects. Judging from the results, the regularizing effect of the Bayesian approach seems to have worked out to at least some degree. From a probabilistic nowcast point of view on the other hand, the skill of BCNN was very lackluster, characterized by a striking lack of sample diversity, unreliability and low discriminatory skill at all precipitation intensity probability thresholds.

Some parts of the model selection and training processes had their issues, leading to a high margin for improvement capacity with possibly little effort, especially in the realm of making reliable uncertainty estimates. Still, these issues are diverse and have unclear roots, meaning that although they might resolve easily, this is not guaranteed.

The silver lining of this work is that the developed models are indeed fast at making inference while making great use of the capacity of Neural Networks to fit nonlinear patterns, and most importantly do so while providing an early iteration of uncertainty estimates. The current model has much potential for improvement as outlined in the Discussion chapter 5. Although it is not yet of big use to forecasters, there is no reason to believe that other more advanced approaches to probabilistic nowcasting using Neu-

ral Networks will not be, whether or not they are based on this work, as there is an ever-increasing demand and need for skillful and reliable probabilistic nowcasting methods.

# Chapter 7

## references

- [1] SVI Part IV: Tips and Tricks - Pyro Tutorials 1.8.1 documentation.
- [2] Next Generation Weather Radar (NEXRAD), Sept. 2020.
- [3] AGRAWAL, S., BARRINGTON, L., BROMBERG, C., BURGE, J., GAZEN, C., AND HICKEY, J. Machine Learning for Precipitation Nowcasting from Radar Images, Dec. 2019. arXiv:1912.12132 [cs, stat].
- [4] ALEXANDER, C., DOWELL, D. C., HU, M., OLSON, J., SMIRNOVA, T., LADWIG, T., WEYGANDT, S., KENYON, J. S., JAMES, E., LIN, H., ET AL. Rapid refresh (rap) and high resolution rapid refresh (hrrr) model development. In *100th American Meteorological Society Annual Meeting* (2020), AMS.
- [5] ANDERSSON, T., AND IVARSSON, K.-I. A model for probability nowcasts of accumulated precipitation using radar. *Journal of Applied Meteorology and Climatology* 30, 1 (1991), 135–141.
- [6] ANDERSSON, T., AND IVARSSON, K.-I. A Model for Probability Nowcasts of Accumulated Precipitation Using Radar. *Journal of Applied Meteorology and Climatology* 30, 1 (Jan. 1991), 135–141. Publisher: American Meteorological Society Section: Journal of Applied Meteorology and Climatology.
- [7] AYZEL, G. Rainnet: a convolutional neural network for radar-based precipitation nowcasting. <https://github.com/hydrogo/rainnet>, 2020.
- [8] AYZEL, G., SCHEFFER, T., AND HEISTERMANN, M. RainNet v1.0: a convolutional neural network for radar-based precipitation nowcasting. 21.

- [9] BAUER, P., THORPE, A., AND BRUNET, G. The quiet revolution of numerical weather prediction. *Nature* 525, 7567 (Sept. 2015), 47–55.
- [10] BINGHAM, E., CHEN, J. P., JANKOWIAK, M., OBERMEYER, F., PRADHAN, N., KARALETOS, T., SINGH, R., SZERLIP, P., HORSFALL, P., AND GOODMAN, N. D. Pyro: Deep Universal Probabilistic Programming. *Journal of Machine Learning Research* (2018).
- [11] BISHOP, C. M., AND NASRABADI, N. M. *Pattern recognition and machine learning*, vol. 4. Springer, 2006.
- [12] BLUNDELL, C., CORNEBISE, J., KAVUKCUOGLU, K., AND WIERSTRAS, D. Weight Uncertainty in Neural Networks. *arXiv:1505.05424 [cs, stat]* (May 2015). arXiv: 1505.05424.
- [13] BOWLER, N. E., PIERCE, C. E., AND SEED, A. W. STEPS: A probabilistic precipitation forecasting scheme which merges an extrapolation nowcast with downscaled NWP. *Quarterly Journal of the Royal Meteorological Society* 132, 620 (2006), 2127–2155. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1256/qj.04.100>.
- [14] BUNTINE, W. Bayesian back-propagation. *Complex systems* 5 (1991), 603–643.
- [15] CAO, Q., KNIGHT, M., FRECH, M., AND MAMMEN, T. Measurement uncertainty and system assessment of weather radar network in Germany. In *2016 IEEE Radar Conference (RadarConf)* (May 2016), pp. 1–5. ISSN: 2375-5318.
- [16] DENKER, J., AND LECUN, Y. Transforming Neural-Net Output Levels to Probability Distributions. In *Advances in Neural Information Processing Systems* (1990), vol. 3, Morgan-Kaufmann.
- [17] DENKER, J., SCHWARTZ, D., WITTNER, B., SOLLA, S., HOWARD, R., JACKEL, L., AND HOPFIELD, J. Large automatic learning, rule extraction, and generalization. *Complex systems* 1, 5 (1987), 877–922.
- [18] DIXON, M., AND WIENER, G. Titan: Thunderstorm identification, tracking, analysis, and nowcasting a radar-based methodology. *Journal of atmospheric and oceanic technology* 10, 6 (1993), 785–797.
- [19] FABRY, F. *Radar Meteorology: Principles and Practice*. Cambridge University Press, Mar. 2018. Google-Books-ID: 0aRwCQAAQBAJ.

- [20] FALCON, W., AND THE PYTORCH LIGHTNING TEAM. PyTorch Lightning, 3 2019.
- [21] FOX, N. I., AND WIKLE, C. K. A bayesian quantitative precipitation nowcast scheme. *Weather and forecasting* 20, 3 (2005), 264–275.
- [22] FRANCH, G., MAGGIO, V., COVIELLO, L., PENDESINI, M., JURMAN, G., AND FURLANELLO, C. TAASRAD19, a high-resolution weather radar reflectivity dataset for precipitation nowcasting. *Scientific Data* 7, 1 (July 2020), 234. Number: 1 Publisher: Nature Publishing Group.
- [23] GAL, Y., AND GHAHRAMANI, Z. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. In *Proceedings of The 33rd International Conference on Machine Learning* (June 2016), PMLR, pp. 1050–1059. ISSN: 1938-7228.
- [24] GERMAN, U., AND ZAWADZKI, I. Scale-dependence of the predictability of precipitation from continental radar images. part i: Description of the methodology. *Monthly Weather Review* 130, 12 (2002), 2859–2873.
- [25] GRAVES, A. Practical Variational Inference for Neural Networks. 9.
- [26] HINTON, E. Keeping Neural Networks Simple by Minimizing the Description Length of the Weights. 9.
- [27] HO, J., KALCHBRENNER, N., WEISSENBORN, D., AND SALIMANS, T. Axial Attention in Multidimensional Transformers, Dec. 2019. arXiv:1912.12180 [cs].
- [28] HOGAN, R. J., FERRO, C. A. T., JOLLIFFE, I. T., AND STEPHENSON, D. B. Equitability Revisited: Why the "Equitable Threat Score" Is Not Equitable. *Weather and Forecasting* 25, 2 (Apr. 2010), 710–726. Publisher: American Meteorological Society Section: Weather and Forecasting.
- [29] KENDALL, A., AND GAL, Y. What Uncertainties Do We Need in Bayesian Deep Learning for Computer Vision? *arXiv:1703.04977 [cs]* (Oct. 2017). arXiv: 1703.04977.
- [30] KINGMA, D. P., SALIMANS, T., AND WELLING, M. Variational Dropout and the Local Reparameterization Trick. *arXiv:1506.02557 [cs, stat]* (Dec. 2015). arXiv: 1506.02557.
- [31] KIUREGHIAN, A. D., AND DITLEVSEN, O. Aleatory or epistemic? Does it matter? *Structural Safety* 31, 2 (Mar. 2009), 105–112.

- [32] KWON, Y., WON, J.-H., KIM, B. J., AND PAIK, M. C. Uncertainty quantification using Bayesian neural networks in classification: Application to biomedical image segmentation. *Computational Statistics & Data Analysis* 142 (Feb. 2020), 106816.
- [33] LUCAS, B. D., KANADE, T., ET AL. *An iterative image registration technique with an application to stereo vision*, vol. 81. Vancouver, 1981.
- [34] MACKAY, D. J. A practical bayesian framework for backpropagation networks. *Neural computation* 4, 3 (1992), 448–472.
- [35] MARSHALL, J., AND PALMER, W. The size distribution of raindrops. *Journal of Atmospheric Sciences* 5 (1948), 165–166.
- [36] MASON, I. A model for assessment of weather forecasts. *Aust. Meteor. Mag* 30, 4 (1982), 291–303.
- [37] MCALLISTER, R., GAL, Y., KENDALL, A., WILK, M. v. D., SHAH, A., CIPOLLA, R., AND WELLER, A. Concrete Problems for Autonomous Vehicle Safety: Advantages of Bayesian Deep Learning. 4745–4753.
- [38] MINKA, T. P. A family of algorithms for approximate Bayesian inference. 75.
- [39] OH, J., GUO, X., LEE, H., LEWIS, R., AND SINGH, S. Action-Conditional Video Prediction using Deep Networks in Atari Games, Dec. 2015. arXiv:1507.08750 [cs].
- [40] OTSUKA, S., TUERHONG, G., KIKUCHI, R., KITANO, Y., TANIGUCHI, Y., RUIZ, J. J., SATOH, S., USHIO, T., AND MIYOSHI, T. Precipitation Nowcasting with Three-Dimensional Space - Time Extrapolation of Dense and Frequent Phased-Array Weather Radar Observations. *Weather and Forecasting* 31, 1 (Feb. 2016), 329–340. Publisher: American Meteorological Society Section: Weather and Forecasting.
- [41] PAN, X., LU, Y., ZHAO, K., HUANG, H., WANG, M., AND CHEN, H. Improving Nowcasting of Convective Development by Incorporating Polarimetric Radar Variables Into a Deep-Learning Model. *Geophysical Research Letters* 48, 21 (2021), e2021GL095302. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1029/2021GL095302>.

- [42] PRUDDEN, R., ADAMS, S., KANGIN, D., ROBINSON, N., RAVURI, S., MOHAMED, S., AND ARRIBAS, A. A review of radar-based nowcasting of precipitation and applicable machine learning techniques. *arXiv:2005.04988 [physics, stat]* (May 2020). arXiv: 2005.04988.
- [43] PULKKINEN, S., CHANDRASEKAR, V., AND NIEMI, T. Lagrangian Integro-Difference Equation Model for Precipitation Nowcasting. *Journal of Atmospheric and Oceanic Technology* 38, 12 (Dec. 2021), 2125–2145. Publisher: American Meteorological Society Section: Journal of Atmospheric and Oceanic Technology.
- [44] PULKKINEN, S., CHANDRASEKAR, V., VON LERBER, A., AND HARRI, A.-M. Nowcasting of Convective Rainfall Using Volumetric Radar Observations. *IEEE Transactions on Geoscience and Remote Sensing* 58, 11 (Nov. 2020), 7845–7859. Conference Name: IEEE Transactions on Geoscience and Remote Sensing.
- [45] PULKKINEN, S., NERINI, D., PÁREZ HORTAL, A. A., VELASCO-FORERO, C., SEED, A., GERMANN, U., AND FORESTI, L. Pysteps: an open-source Python library for probabilistic precipitation nowcasting (v1.0). *Geoscientific Model Development* 12, 10 (Oct. 2019), 4185–4219. Publisher: Copernicus GmbH.
- [46] RAVURI, S., LENCI, K., WILLSON, M., KANGIN, D., LAM, R., MIROWSKI, P., FITZSIMONS, M., ATHANASSIADOU, M., KASHEM, S., MADGE, S., PRUDDEN, R., MANDHANE, A., CLARK, A., BROCK, A., SIMONYAN, K., HADSELL, R., ROBINSON, N., CLANCY, E., ARRIBAS, A., AND MOHAMED, S. Skilful precipitation nowcasting using deep generative models of radar. *Nature* 597, 7878 (Sept. 2021), 672–677. Number: 7878 Publisher: Nature Publishing Group.
- [47] RINEHART, R. E., AND GARVEY, E. T. Three-dimensional storm motion detection by conventional weather radar. *Nature* 273, 5660 (May 1978), 287–289. Number: 5660 Publisher: Nature Publishing Group.
- [48] RITTER, H., AND KARALETOS, T. TyXe: Pyro-based Bayesian neural nets for Pytorch. *arXiv preprint arXiv:2110.00276* (2021).
- [49] ROBERTS, N. M., AND LEAN, H. W. Scale-Selective Verification of Rainfall Accumulations from High-Resolution Forecasts of Convective Events. *Monthly Weather Review* 136, 1 (Jan. 2008), 78–97. Publisher: American Meteorological Society Section: Monthly Weather Review.

- [50] RYU, S., LYU, G., DO, Y., AND LEE, G. Improved rainfall nowcasting using Burgers' equation. *Journal of Hydrology* 581 (Feb. 2020), 124140.
- [51] SAKAINO, H. Spatio-Temporal Image Pattern Prediction Method Based on a Physical Model With Time-Varying Optical Flow. *IEEE Transactions on Geoscience and Remote Sensing* 51, 5 (May 2013), 3023–3036. Conference Name: IEEE Transactions on Geoscience and Remote Sensing.
- [52] SALTIKOFF, E., HAASE, G., DELOBBE, L., GAUSSIAT, N., MARTET, M., IDZIOREK, D., LEIJNSE, H., NOVÁJK, P., LUKACH, M., AND STEPHAN, K. OPERA the Radar Project. *Atmosphere* 10, 6 (June 2019), 320. Number: 6 Publisher: Multidisciplinary Digital Publishing Institute.
- [53] SCHAEFER, J. T. The Critical Success Index as an Indicator of Warning Skill. *Weather and Forecasting* 5, 4 (Dec. 1990), 570–575. Publisher: American Meteorological Society Section: Weather and Forecasting.
- [54] SCHMID, F., WANG, Y., AND HAROU, A. Nowcasting guidelines—a summary. *Bulletin n°* 68 (2019), 2.
- [55] SCHMID, W., MECKLENBURG, S., AND JOSS, J. Short-term risk forecasts of heavy rainfall. *Water science and technology* 45, 2 (2002), 121–125.
- [56] SCHULTZ, M. G., BETANCOURT, C., GONG, B., KLEINERT, F., LANGGUTH, M., LEUFEN, L. H., MOZAFFARI, A., AND STADTLER, S. Can deep learning beat numerical weather prediction? *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 379, 2194 (Apr. 2021), 20200097. Publisher: Royal Society.
- [57] SEED, A. W. A Dynamic and Spatial Scaling Approach to Advection Forecasting. *Journal of Applied Meteorology and Climatology* 42, 3 (Mar. 2003), 381–388. Publisher: American Meteorological Society Section: Journal of Applied Meteorology and Climatology.
- [58] SEED, A. W., PIERCE, C. E., AND NORMAN, K. Formulation and evaluation of a scale decomposition-based stochastic precipitation nowcast scheme. *Water Resources Research* 49, 10 (2013), 6624–6641. eprint: <https://agupubs.onlinelibrary.wiley.com/doi/pdf/10.1002/wrcr.20536>.

- [59] SHI, X., CHEN, Z., WANG, H., YEUNG, D.-Y., WONG, W.-K., AND WOO, W.-C. Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting, Sept. 2015. arXiv:1506.04214 [cs] version: 2.
- [60] SHI, X., GAO, Z., LAUSEN, L., WANG, H., YEUNG, D.-Y., WONG, W.-K., AND WOO, W.-C. Deep Learning for Precipitation Nowcasting: A Benchmark and A New Model, Oct. 2017. arXiv:1706.03458 [cs].
- [61] SHRIDHAR, K., LAUMANN, F., AND LIWICKI, M. A Comprehensive guide to Bayesian Convolutional Neural Network with Variational Inference. *arXiv:1901.02731 [cs, stat]* (Jan. 2019). arXiv: 1901.02731.
- [62] SÄNDERBY, C. K., ESPEHOLT, L., HEEK, J., DEHGHANI, M., OLIVER, A., SALIMANS, T., AGRAWAL, S., HICKEY, J., AND KALCHBRENNER, N. MetNet: A Neural Weather Model for Precipitation Forecasting, Mar. 2020. arXiv:2003.12140 [physics, stat].
- [63] TISHBY, LEVIN, AND SOLLA. Consistent inference of probabilities in layered networks: predictions and generalizations. In *International 1989 Joint Conference on Neural Networks* (1989), pp. 403–409 vol.2.
- [64] VOORMANSIK, T., ROSSI, P. J., MOISSEEV, D., TANILSOO, T., AND POST, P. Thunderstorm hail and lightning detection parameters based on dual-polarization Doppler weather radar data. *Meteorological Applications* 24, 3 (2017), 521–530. \_eprint: <https://rmets.onlinelibrary.wiley.com/doi/pdf/10.1002/met.1652>.
- [65] WANG, Z., SIMONCELLI, E., AND BOVIK, A. Multiscale structural similarity for image quality assessment. In *The Thrity-Seventh Asilomar Conference on Signals, Systems & Computers, 2003* (Pacific Grove, CA, USA, 2003), IEEE, pp. 1398–1402.
- [66] WEN, Y., VICOL, P., BA, J., TRAN, D., AND GROSSE, R. Flipout: Efficient Pseudo-Independent Weight Perturbations on Mini-Batches, Apr. 2018. Number: arXiv:1803.04386 arXiv:1803.04386 [cs, stat].
- [67] YIN, J., GAO, Z., AND HAN, W. Application of a Radar Echo Extrapolation-Based Deep Learning Method in Strong Convection Nowcasting. *Earth and Space Science* 8, 8 (2021), e2020EA001621. \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1029/2020EA001621>.
- [68] ZHAO, H., GALLO, O., FROSIO, I., AND KAUTZ, J. Loss Functions for Image Restoration With Neural Networks. *IEEE Transactions on*

*Computational Imaging* 3, 1 (Mar. 2017), 47–57. Conference Name: IEEE Transactions on Computational Imaging.



# Appendix A

## Additional Figures

### A.1 Model Selection

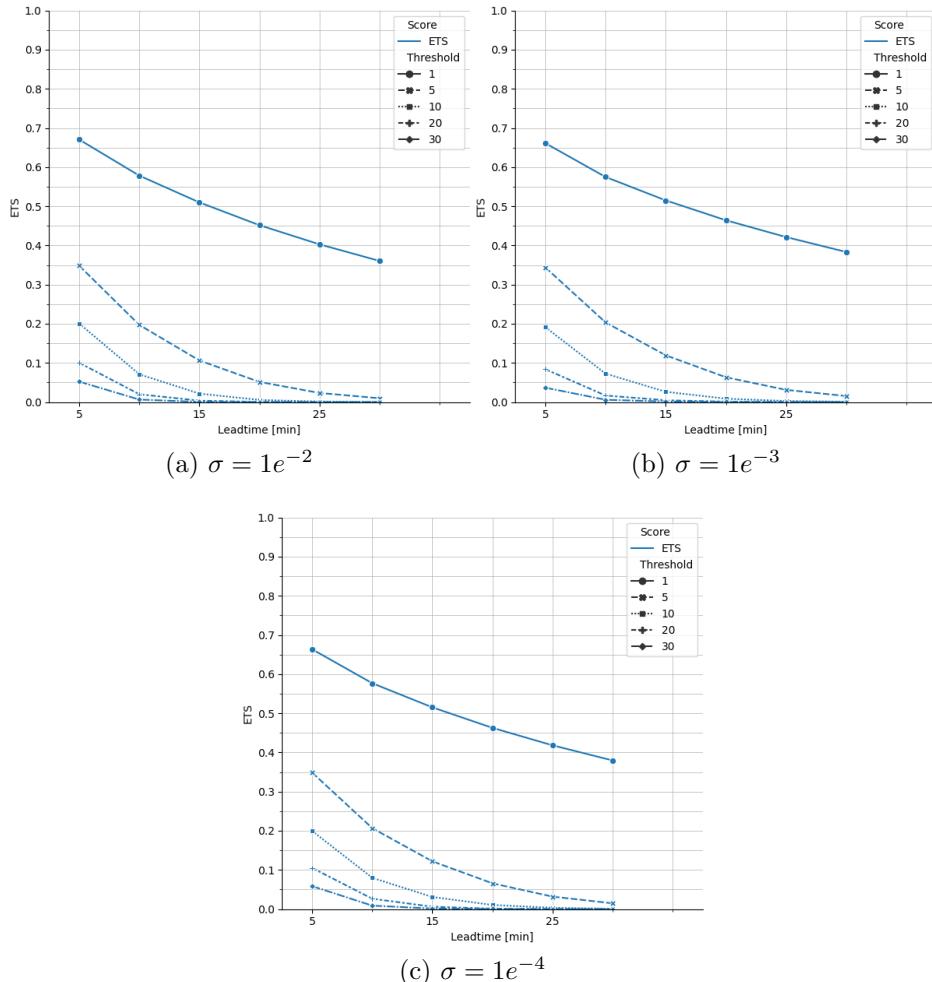


Figure A.1: ETS scores of RainNet over the validation set at model convergence used for choosing a  $\sigma$  Gaussian NLL parameter. Thresholds are expressed in mm/h.

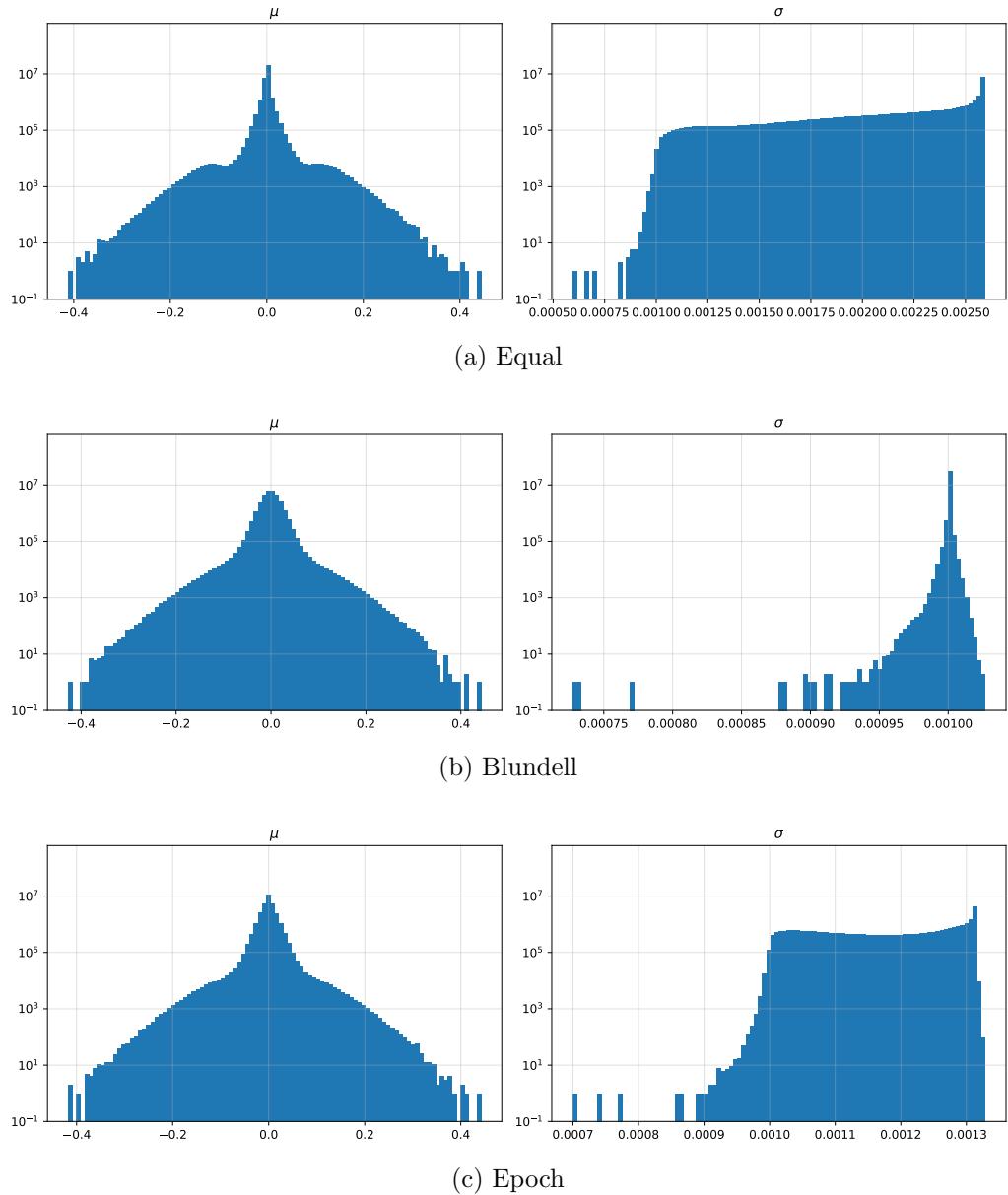


Figure A.2: The relation of KL-weighting scheme to the parameter posterior distributions across the network.  $\mu$  histograms describe the parameter mean distributions, whereas  $\sigma$  distributions describe the distribution of their standard deviations. The prior chosen is GSM SD.

## A.2 Case Studies

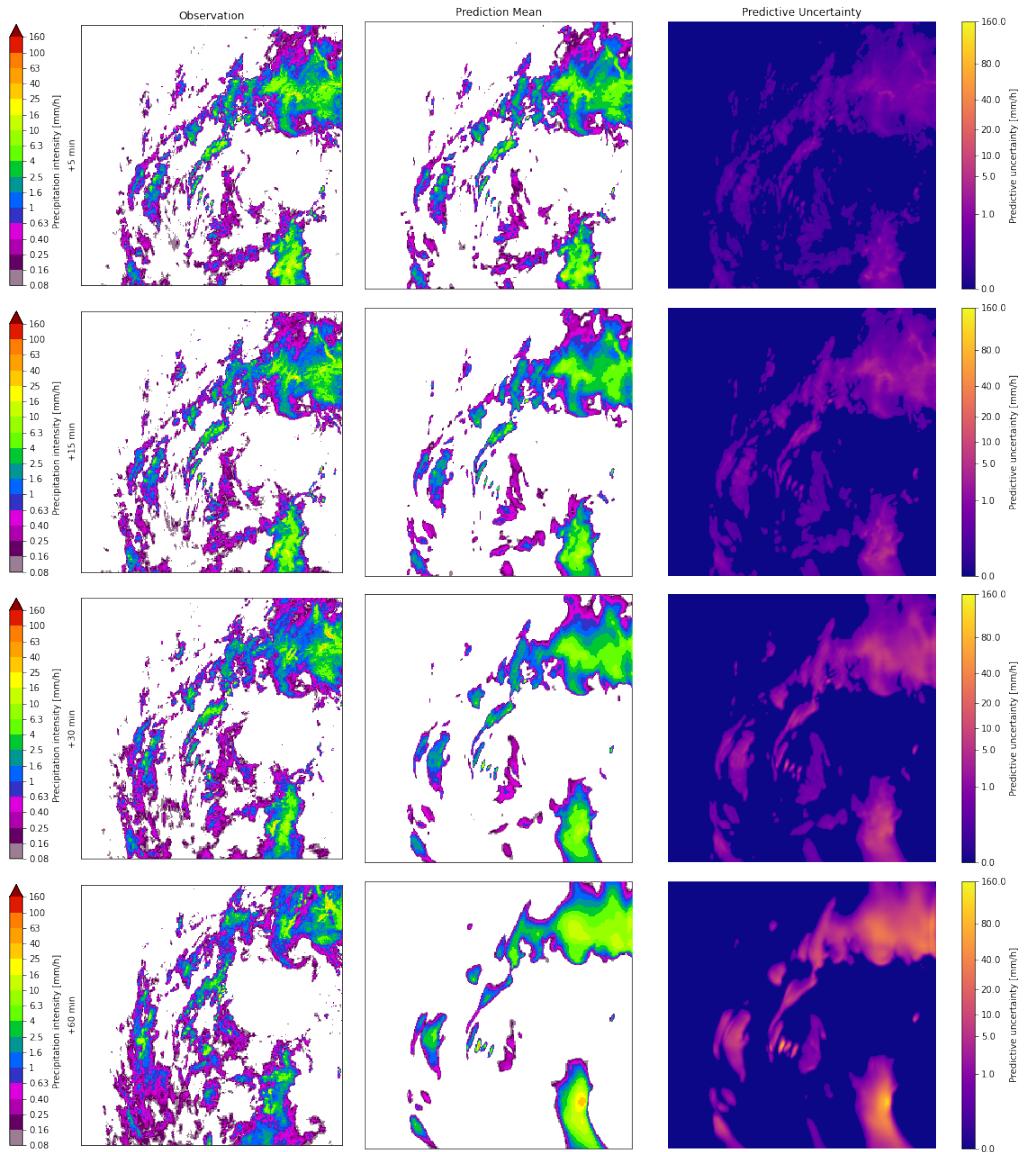


Figure A.3: Predictive mean and uncertainty (two standard deviations) for Case 1 on the 25th of May 2019 starting at 1PM for the BCNN 1t5 new model.

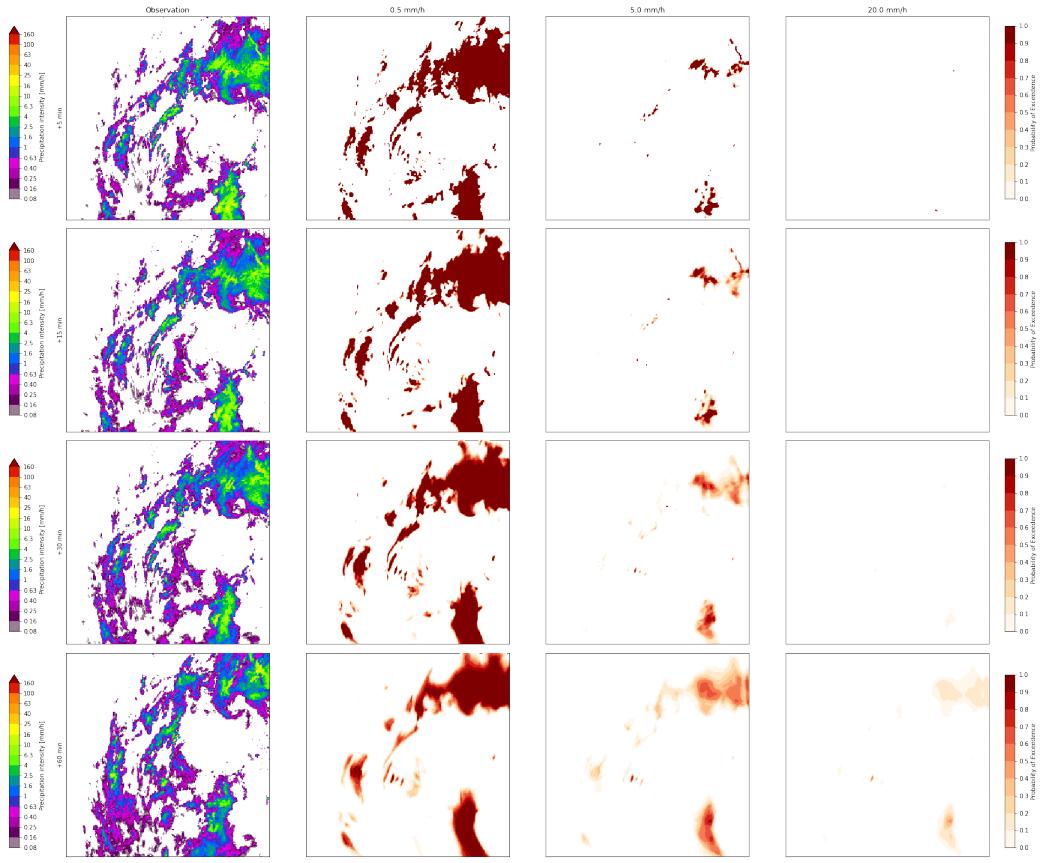


Figure A.4: 0.5, 5.0, and 20.0 mm/h exceedance probabilities for Case 1 on the 25th of May 2019 starting at 1PM for the BCNN 1t5 new model.

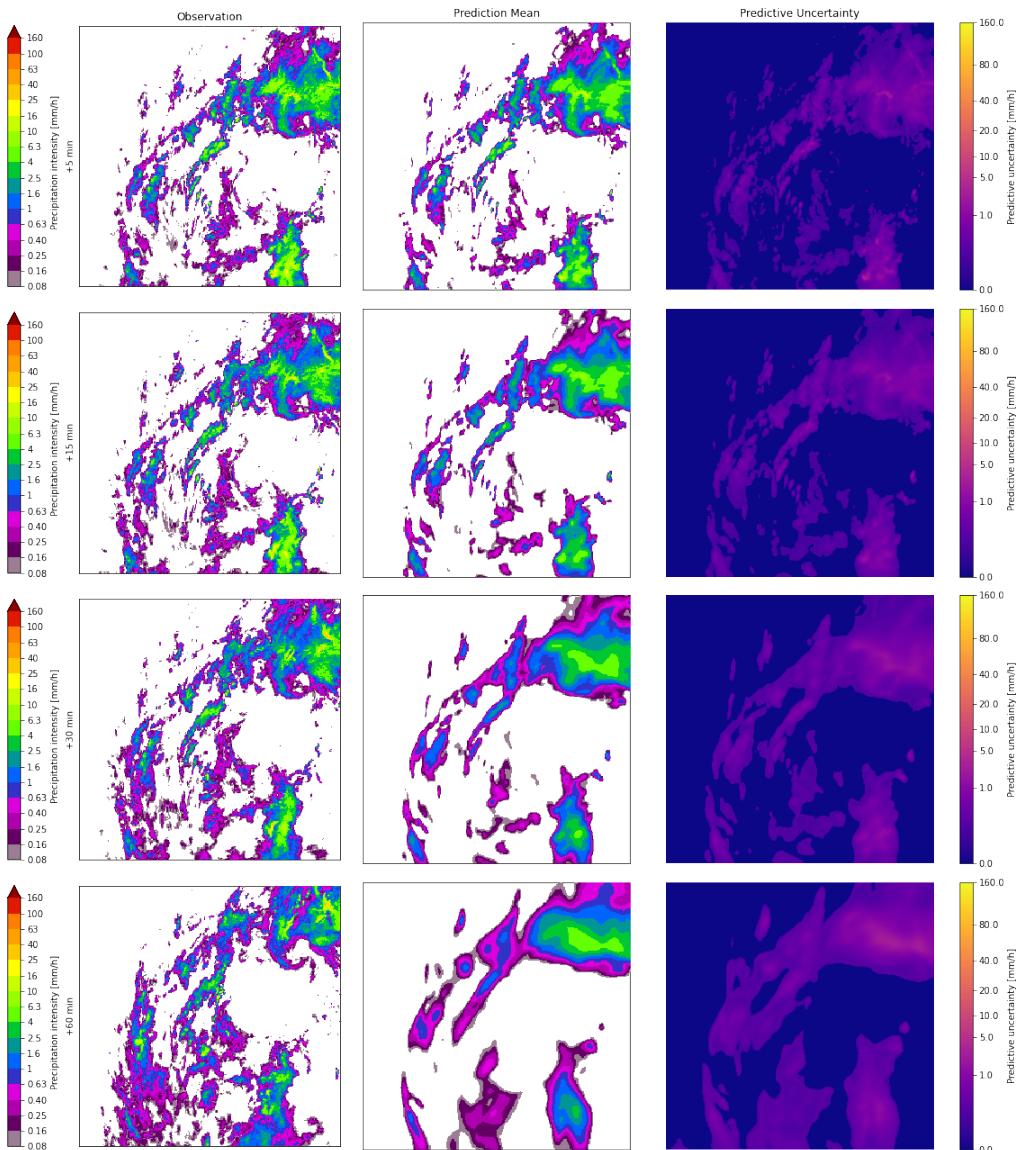


Figure A.5: Predictive mean and uncertainty (two standard deviations) for Case 1 on the 25th of May 2019 starting at 1PM for the BCNN 1t30 model.

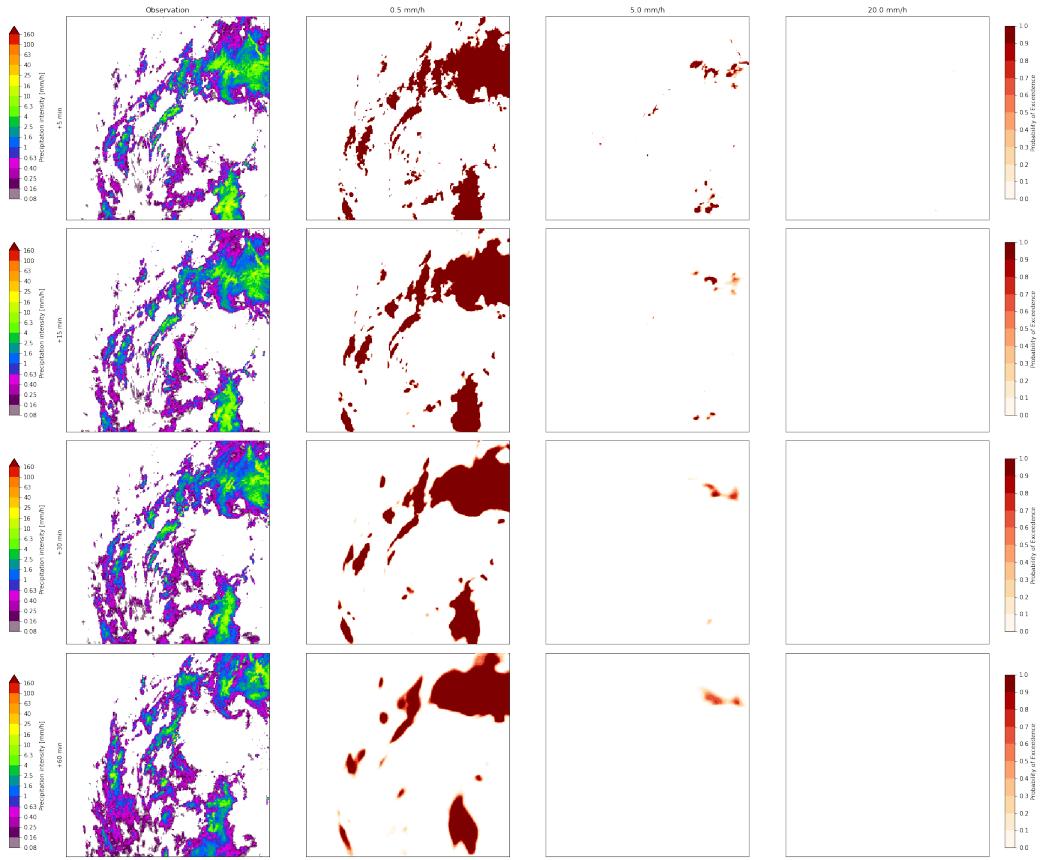


Figure A.6: 0.5, 5.0, and 20.0 mm/h exceedance probabilities for Case 1 on the 25th of May 2019 starting at 1PM for the BCNN 1t30 model.

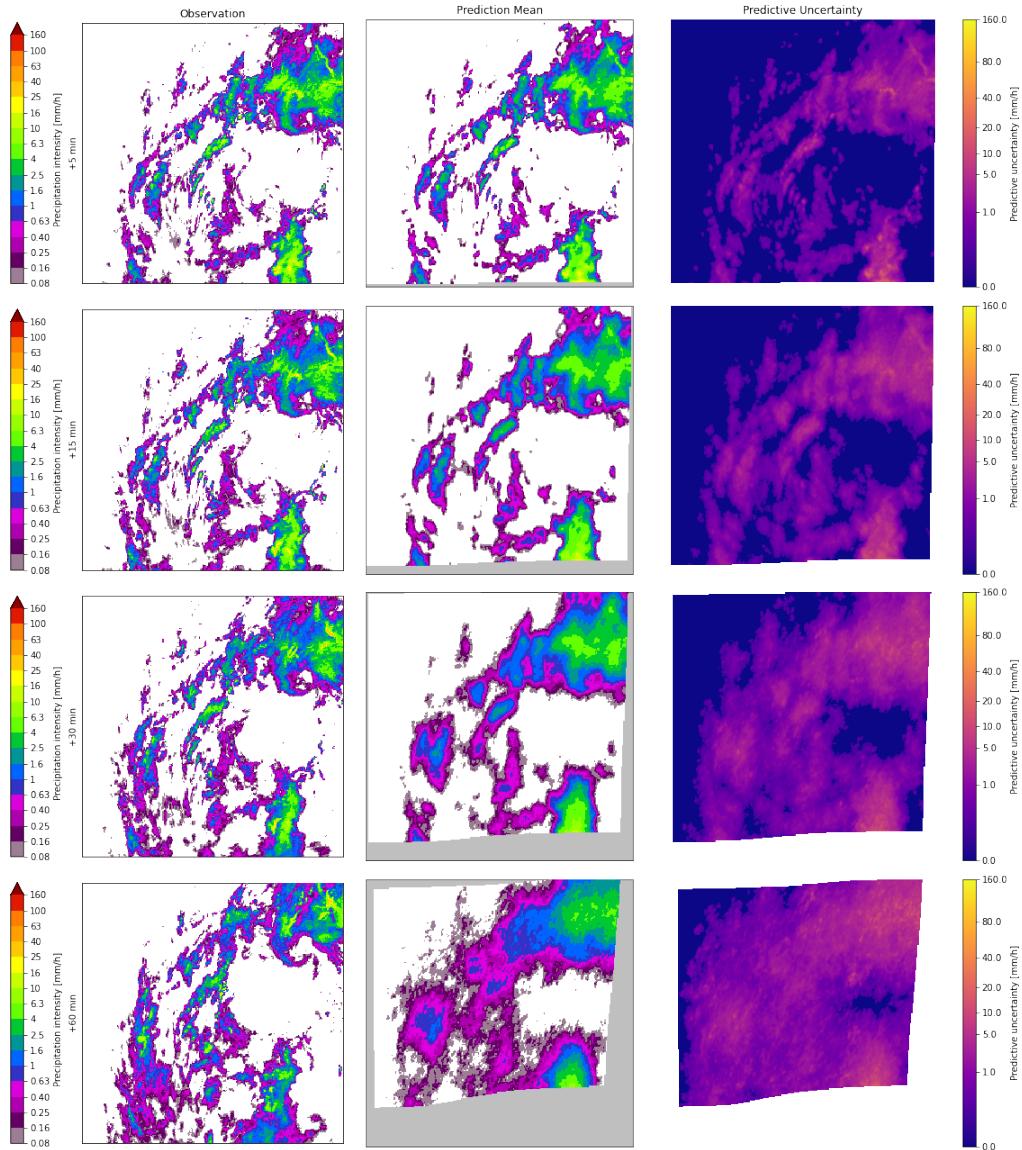


Figure A.7: Predictive mean and uncertainty (two standard deviations) for Case 1 on the 25th of May 2019 starting at 1PM for the baseline STEPS model.

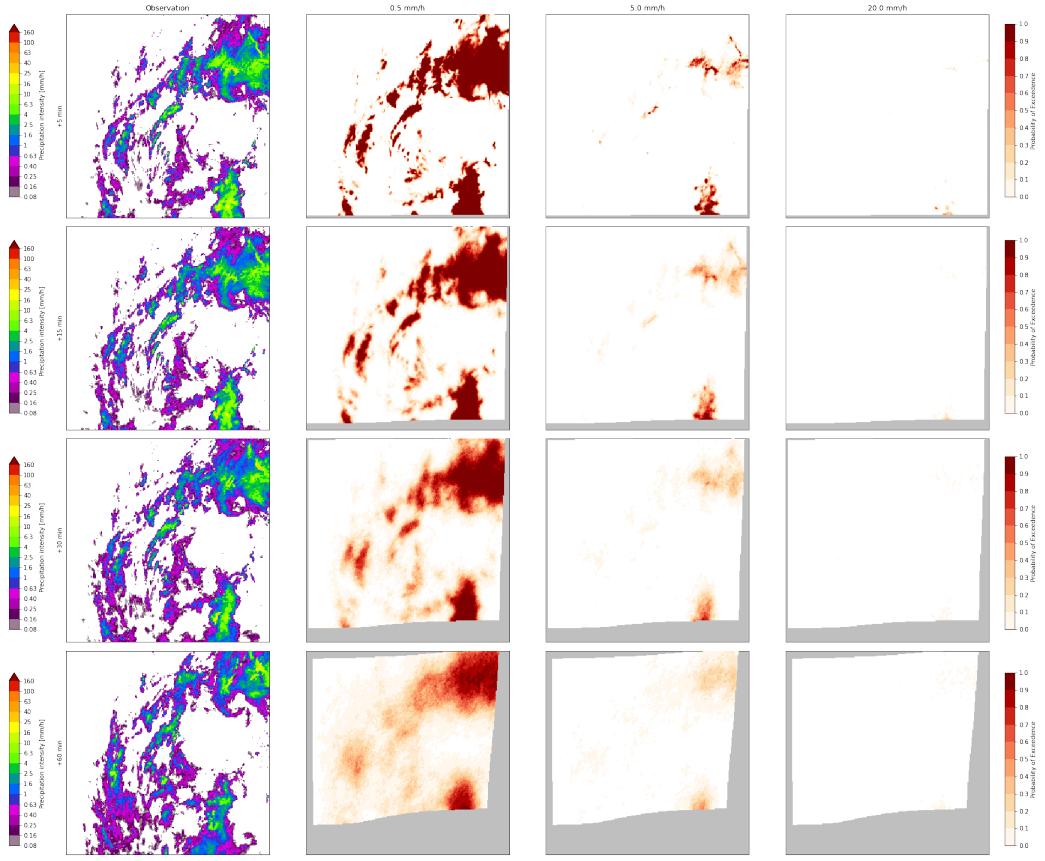


Figure A.8: 0.5, 5.0, and 20.0 mm/h exceedance probabilities for Case 1 on the 25th of May 2019 starting at 1PM for the baseline STEPS model.

### A.3 Deterministic Metrics

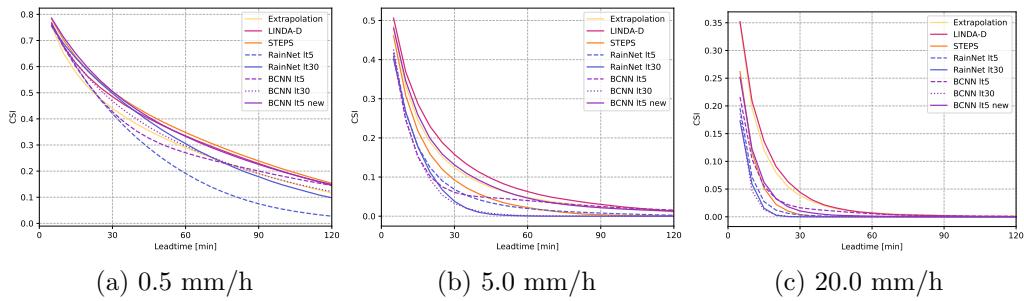


Figure A.9: CSI scores

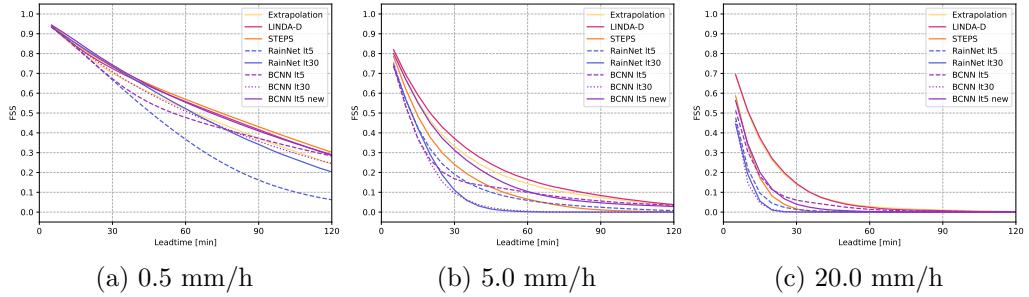


Figure A.10: FSS scores (4km)

## A.4 Probabilistic Metrics

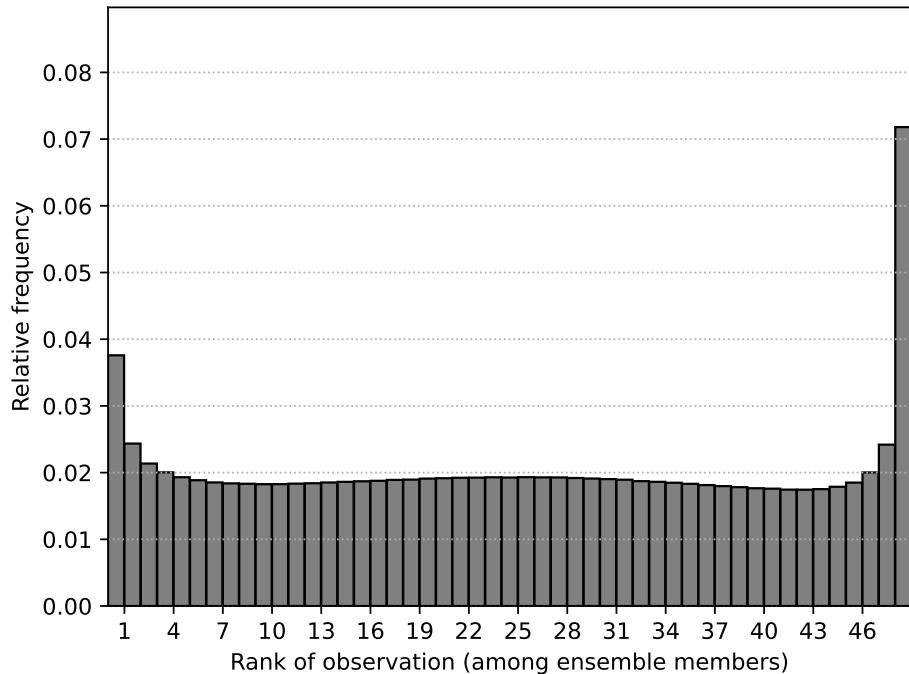


Figure A.11: Rank Histogram for LINDA-P model nowcasts at a one-hour leadtime, for precipitation events exceeding 0.1 mm/h