

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! — June 22, 2022 — DRAFT!**

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

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

ABSTRACT OF  
MASTER'S THESIS

<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>	33
<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.		
!FIXME This is an example how to use fixme: add your abstract here. FIXME!			
<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 Neuro-verkkoja		
<b>Päiväys:</b>	20. heinäkuuta 2022	<b>Sivumäärä:</b>	33
<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.		
joku			
<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
BRN	Bayesian RainNet
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	Equivalent 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

# 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 Weather radars and radar products . . . . .	11
2.1.2 Overview of weather forecast methods . . . . .	11
2.1.3 Precipitation nowcasting : classical deterministic meth-	
ods . . . . .	12
2.1.4 Precipitation nowcasting : classical probabilistic methods	12
2.1.5 Machine learning approaches to precipitation nowcasting	12
2.2 Bayesian deep learning . . . . .	12
2.2.1 Learning probability distributions . . . . .	12
2.2.2 Variational inference . . . . .	12
2.2.3 Predictive uncertainty estimation and decomposition .	14
2.3 Related work . . . . .	15
<b>3 Materials and Methods</b>	<b>16</b>
3.1 Dataset and data selection . . . . .	16
3.2 Model . . . . .	17
3.2.1 The baseline: RainNet . . . . .	17
3.2.2 Our model: a bayesian extension to RainNet . . . . .	17
3.3 Verification Methods . . . . .	18
3.3.1 Evaluation of nowcast predictive uncertainty . . . . .	18
3.3.2 Baseline Deterministic Models . . . . .	18
3.3.3 Baseline Probabilistic Models . . . . .	18
3.3.4 Deterministic Prediction skill evaluation metrics . . . .	18
3.3.5 Probabilistic Prediction skill evaluation metrics . . . .	21

3.3.6	Practical points regarding verification experiments . . .	23
3.4	Software and resources . . . . .	23
<b>4</b>	<b>Results</b>	<b>25</b>
4.1	Case studies for nowcasts . . . . .	25
4.1.1	Case study 1 : Large-scale Stratiform rain event . . . .	25
4.1.2	Case study 2 : Rapidly evolving convective rain event .	25
4.2	Deterministic prediction skill (Metrics) . . . . .	25
4.3	Probabilistic prediction skill (Metrics) . . . . .	25
4.4	Uncertainty estimation . . . . .	25
4.4.1	Uncertainties against leadtime . . . . .	25
4.4.2	Uncertainties against rainfall intensity . . . . .	25
4.4.3	Epistemic uncertainty against training parameters . . .	25
<b>5</b>	<b>Discussion</b>	<b>26</b>
5.1	Goodness of results . . . . .	26
5.2	Validity of results . . . . .	26
5.3	What could we learn from uncertainty . . . . .	26
5.4	What would have to be improved, potential problems in the study? . . . . .	26
5.5	Directions for further work . . . . .	26
<b>6</b>	<b>Conclusions</b>	<b>28</b>
<b>7</b>	<b>references</b>	<b>29</b>
<b>A</b>	<b>First appendix</b>	<b>32</b>

# Chapter 1

## Introduction

**!FIXME Add rest of references, by 22.6, 45min FIXME!**

**!FIXME remove "we" passive FIXME!**

Nowcasting is defined as weather forecasting at very short timescales, from minutes ahead up to several hours. 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, flash flood prevention, and mitigation of economic loss incurred by sudden and heavy precipitation. 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. 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 their computational expensiveness and the low spatial resolution of large-scale NWP, making those systems unfit to predict for example heavy localized rainfall.

**!FIXME add a tiny little bit about det classical nowcast methods as tongue teaser, by 22.6, 45min FIXME!**

To compensate for the shortcomings of NWP in the realm of precipitation nowcasting, many different dedicated models and algorithms have been developed. 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 advection field. Other methods have been developed that track and try to nowcast the evolution of indi-



vidual rain cells, instead of the whole field. 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 convective initiation and life-cycles. 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, their expressive power and the relative cheapness of inference.

## 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 phenomena we are most interested in are extreme events, and these are the events that we want to prepare for, it makes intuitive sense that accurate and reliable probabilistic nowcasts are primordial for operational decision-making in meteorological crises.

Indeed, the interest in probabilistic nowcasts has grown lately and many new probabilistic models have been introduced in the last decades, notably the Short-Term Ensemble Prediction System (STEPS) [5] and more recently the Lagrangian Integro-Difference equation model with Autoregression (LINDA) [10]. These methods predict ensembles, that are used to estimate underlying distributions of data and precipitation probabilities. They perform reasonably well, but have computationally expensive inference and suffer from the same limitations other advection-extrapolation based methods.

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. [12] 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 we will focus on from now on is instead to model the uncertainty of model parameters rather than that of the data. 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 [7] and autonomous vehicles [8].

In this work, we will approach the problem of making skillful and reliable probabilistic precipitation nowcasts by building an uncertainty-aware Bayesian Neural Networks. A baseline Convolutional Neural Network (CNN) will be turned into a Bayesian Convolutional Network (BCN) with optimization performed using stochastic variational inference (SVI) over model parameters. The model is trained, validated and tested using Finnish Meteorological Institute (FMI) radar reflectivity composites. For 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

!FIXME Estimate at start: 12-15 pages FIXME!

## 2.1 Precipitation Nowcasting

### 2.1.1 Weather radars and radar products

- physical functioning and paragraph on history of weather radar
- reflectivity products, 2D and 3D
- from reflectivity to RR, limitations of reflectivity
- double-polarization weather radars and Polarimetric products: info on microphysical processes

$$z = AR^b \tag{2.1}$$

where  $z$  is reflectivity,  $R$  is rainrate, and  $A$  as well as  $b$  are empirically determined parameters. Reflectivity is calculated

### 2.1.2 Overview of weather forecast methods

- Taking a broad look into the NWP process: Data collection, data assimilation, numerical prediction, postprocessing into forecast products
- NWP in the context of precipitation forecasts

### **2.1.3 Precipitation nowcasting : classical deterministic methods**

- Basic principles of advection equations.
- Chronological advancing, with +/- of each method

### **2.1.4 Precipitation nowcasting : classical probabilistic methods**

- approach types, ensemble etc

### **2.1.5 Machine learning approaches to precipitation nowcasting**

- 

## **2.2 Bayesian deep learning**

### **2.2.1 Learning probability distributions**

- What motivation, first old studies
- overview on Intractable integrals and ways to deal with them
- MCMC, etc...
- finish with switch to VI

### **2.2.2 Variational inference**

- deriving ELBO
- Bayes by Backprop and SVI
- Local reparametrization
- Monte-Carlo dropout as VI
- where else is VI used
- problems with VI

The evidence lower bound (ELBO), also known as variational free energy, is thus defined as

$$\mathcal{F}(\mathcal{D}, \sigma) = \text{KL}[q(\mathbf{w}|\sigma)||P(\mathbf{w})] - \mathbb{E}_{q(\mathbf{w}|\sigma)}[\log(P(\mathcal{D}|\mathbf{w}))] \quad (2.2)$$

In order to apply backpropagation, the standard deviation is parametrized 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 allows gradients to flow through the network as sampling operations are external.

- how does variance scale for basic re-parametrization trick?
- more Details on LRT

Kingma et al. (2015) [?] introduced a modification to the reparametrization trick, called the *local reparametrization trick* (LRT). It works similarly but layer activations are sampled using  $\epsilon$  rather than weights in order to reduce the computational cost, and even more importantly variances in a minibatch, allowing for faster training.

- Flipout and its possible advantages over LRT

Using the reparametrization trick and

$$\frac{\partial}{\partial \sigma} \mathbb{E}_{q(\mathbf{w}|\sigma)}[f(\mathbf{w}, \sigma)] = \mathbb{E}_{q(\epsilon)} \left[ \frac{\partial f(\mathbf{w}, \sigma)}{\partial \mathbf{w}} + \frac{\partial \mathbf{w}}{\partial \sigma} \frac{\partial f(\mathbf{w}, \sigma)}{\partial \mathbf{w}} \right] \quad (2.3)$$

the ELBO cost function (Eq. 2.2) can be rewritten as in a form amenable to minibatch optimization, yielding

- Elaborate on the above proposition by Blundell et al

$$\mathcal{F}(\mathcal{D}, \sigma) \approx \sum_{i=1}^n \log q(\mathbf{w}^{(i)}|\sigma) - \log P(\mathbf{w}^{(i)}) - \log P(\mathcal{D}|\mathbf{w}^{(i)}) \quad (2.4)$$

Here  $\mathbf{w}^{(i)}$  denotes the  $i$ :th monte-carlo sample drawn the the posterior. From the terms in eq. 2.4:

1.  $\log q(\mathbf{w}^{(i)}|\sigma)$  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 also often called the complexity cost.

The cost to be minimizing can be again rewritten as

$$\mathcal{F}_i^\pi(\mathcal{D}, \sigma) = \pi_i \text{KL}[q(\mathbf{w}|\sigma)||P(\mathbf{w})] - \mathbb{E}_{q(\mathbf{w}|\sigma)}[\log P(\mathcal{D}_i|\mathbf{w})] \quad (2.5)$$

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) [?] used equal weighting as  $\pi_i = 1/M$ , but Blundell et al. (2015) [?] 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 [16?] because of their properties facilitating optimization. These scale mixture priors, if containing two scales are defined as

$$P(\mathbf{w}) = \prod_j \alpha \mathcal{N}(\mathbf{w}_j|0, \sigma_1^2) + (1 - \alpha) \mathcal{N}(\mathbf{w}_j|0, \sigma_2^2) \quad (2.6)$$

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.

### 2.2.3 Predictive uncertainty estimation and decomposition

!FIXME find good ref on sources !FIXME!

- Where does uncertainty come from
- division into epistemic, aleatoric uncertainty in ML
- Division for classification
- Division for regression

## 2.3 Related work

`!FIXME` Section to include if 2.1.4 and 2.1.5 become convoluted  
`FIXME!`

## 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 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.

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 made of these radar scans are readily available at FMI which facilitated their retrieval for this work. Additionally, the bounding box covering southern Finland did not have any major data quality issues standing out as opposed to other candidate datasets, and finally missing pixels were correctly labeled. All of these factors contributed to the choice of the data source.

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.

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 [15] with 6 hour long blocks to prevent auto-correlation between consecutive radar images from invalidating independence 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.

Radar composites are stored as gzip-compressed PGM composite files holding uint8 values, which are converted to reflectivity (dBZ) by applying the relation  $\text{pixel\_value}_{\text{dBZ}} = -32 + 0.5 * \text{pixel\_value}_{\text{uint8}}$ . For being fed into neural networks, the composites are read from HDF5 files as this improves reading speed on distributed storage systems.

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

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

!FIXME Add domain presenting image ?    !FIXME!

## 3.2 Model

### 3.2.1 The baseline: RainNet

For the implementation of a Bayesian Convolutional Network, we use as our base model RainNet [2], which is itself heavily inspired by UNet and SegNet model families.

!FIXME Add RainNet training and/or architecture diagrams    !FIXME!

- Describe loss functions and training procedure
- describe mods from UNet, SegNet, its problems, good points
- describe modifications made for this work

### 3.2.2 Our model: a bayesian extension to RainNet

- Converting weights into distributions
- Gaussian prior, Gaussian posterior enabling closed-form KL-divergence
- Data likelihood modeling: Homeoscedastic Gaussian likelihood
- Training procedure and the Local Reparametrization Trick used.

### 3.3 Verification Methods

#### 3.3.1 Evaluation of nowcast predictive uncertainty

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.

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 24 for preliminary small scale models and 48 for full-scale models. Verification might be even 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.

#### 3.3.2 Baseline Deterministic Models

- Lagrangian persistence (advection-extrapolation)
- S-PROG
- ANVIL
- LINDA-D
- RainNet (c.f. above)

#### 3.3.3 Baseline Probabilistic Models

- STEPS
- LINDA-P

#### 3.3.4 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 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

$$\text{ME} = \frac{1}{N} \sum_{i=1}^N y_i - \hat{y}_i, \quad (3.2)$$

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,

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|, \quad (3.3)$$

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. 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.

**!FIXME small picture showing contingency table to help visualization** FIXME!

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. Scores derived from those quantities that were used are probability of detection (POD) defined as

$$\text{POD} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (3.4)$$

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

$$\text{FAR} = \frac{\text{FP}}{\text{TP} + \text{FP}} \quad (3.5)$$

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

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

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) defined as

$$\begin{aligned} \text{ETS} &= \frac{\text{TP} - \text{rnd}}{\text{TP} + \text{FN} + \text{FP} - \text{rnd}}, \\ \text{where } \text{rnd} &= \frac{(\text{TP} + \text{FN})(\text{TP} + \text{FP})}{\text{TP} + \text{FN} + \text{FP} + \text{TN}} \end{aligned} \quad (3.7)$$

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 *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) 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 [14]. 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 results 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!**

### 3.3.5 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 breath 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

$$\text{CRPS}(F, x) = \int_{-\infty}^{\infty} (F(y) - \mathbf{1}(x \geq y))^2 dy \quad (3.8)$$

where  $F$  is the forecast cumulative density function (CDF) and  $\mathbf{1}(x \geq y)$

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, although ME tells the same thing.

ROC curves 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 forecast probabilities. When these two values are close to each others, the forecast is said to be reliable. This means that a skillful forecast is defined by being close to the line of equality to observed relative frequencies for all forecast probabilities output by the model.

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!**

### 3.3.6 Practical points regarding verification experiments

All of the scores described in the previous sections 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.

In practice, predictions are calculated for all models for all the timestamps contained in the test set described in section 3.1. Predictions are computed for 36 timesteps, which is equivalent to 3 hours into the future with an interval of 5 minutes. After that, verification metrics are computed using the predictions of each model and they 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.

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.

## 3.4 Software and resources

**!FIXME Some diagram for visualizing grand scheme of workflow  
FIXME!**

The deep learning models were all implemented using the PyTorch framework and the PyTorch Lightning wrapper [6]. 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. [1].

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 [3], which is itself built on top of PyTorch and includes a fully-featured implementation of Stochastic Variational Inference (SVI). In order to facilitate the implementation task, the TyXe package [13] 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 [11]. 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.



## Chapter 4

# Results

### 4.1 Case studies for nowcasts

#### 4.1.1 Case study 1 : Large-scale Stratiform rain event

#### 4.1.2 Case study 2 : Rapidly evolving convective rain event

### 4.2 Deterministic prediction skill (Metrics)

### 4.3 Probabilistic prediction skill (Metrics)

### 4.4 Uncertainty estimation

#### 4.4.1 Uncertainties against leadtime

#### 4.4.2 Uncertainties against rainfall intensity

#### 4.4.3 Epistemic uncertainty against training parameters

## Chapter 5

# Discussion

### 5.1 Goodness of results

### 5.2 Validity of results

Variational approximations tend to underestimate uncertainty of learned distributions. [4, 9]. This can be observed from the rank histograms of BRN as a generally concave histogram shape compared to models based on stochastic perturbations. This means that the effect of variability of ensembles being less than that of observations was most indeed most pronounced in variational models, as one might expect.

### 5.3 What could we learn from uncertainty

### 5.4 What would have to be improved, potential problems in the study?

### 5.5 Directions for further work

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.

## Chapter 6

# Conclusions

## Chapter 7

## references

- [1] AYZEL, G. Rainnet: a convolutional neural network for radar-based precipitation nowcasting. <https://github.com/hydrogo/rainnet>, 2020.
- [2] AYZEL, G., SCHEFFER, T., AND HEISTERMANN, M. RainNet v1.0: a convolutional neural network for radar-based precipitation nowcasting. 21.
- [3] BINGHAM, E., CHEN, J. P., JANKOWIAK, M., OBERMEYER, F., PRADHAN, N., KARALETSOS, T., SINGH, R., SZERLIP, P., HORSFALL, P., AND GOODMAN, N. D. Pyro: Deep Universal Probabilistic Programming. *Journal of Machine Learning Research* (2018).
- [4] BISHOP, C. M., AND NASRABADI, N. M. *Pattern recognition and machine learning*, vol. 4. Springer, 2006.
- [5] 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>.
- [6] FALCON, W., AND THE PYTORCH LIGHTNING TEAM. PyTorch Lightning, 3 2019.
- [7] 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.

- [8] 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.
- [9] MINKA, T. P. A family of algorithms for approximate Bayesian inference. 75.
- [10] 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.
- [11] PULKKINEN, S., NERINI, D., PÁ©REZ 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.
- [12] RAVURI, S., LENC, 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.
- [13] RITTER, H., AND KARALETOS, T. TyXe: Pyro-based Bayesian neural nets for Pytorch. *arXiv preprint arXiv:2110.00276* (2021).
- [14] 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.
- [15] 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.

- [16] 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.

## Appendix A

### First appendix

This is the first appendix. You could put some test images or verbose data in an appendix, if there is too much data to fit in the actual text nicely.

For now, the Aalto logo variants are shown in Figure A.1.





(a) In English



(b) Suomeksi



(c) På svenska

Figure A.1: Aalto logo variants