

1 **A calibrated and consistent combination of probabilistic forecasts**
2 **for the upcrossing of several precipitation thresholds using neural**
3 **networks**

4 P. Schaumann¹, R. Hess², M. Rempel², U. Blahak² and V. Schmidt¹

5 ¹*Institute of Stochastics, Ulm University, Ulm, Germany*

6 ²*Deutscher Wetterdienst, Offenbach, Germany*

ABSTRACT

8 The seamless combination of nowcasting and numerical weather prediction (NWP) aims to
9 provide a functional basis for very-short-term forecasts, that are essential e.g. for weather
10 warnings. In this paper we propose a statistical method for precipitation using neural net-
11 works (NN) that combines nowcasting data from DWD's radar based RadVOR system with
12 postprocessed forecasts of the high resolving NWP ensemble COSMO-DE-EPS. The post-
13 processing is performed by Ensemble-MOS of DWD. Whereas the quality of the nowcasting
14 projections of RadVOR is excellent at the beginning, it declines rapidly after about 2 hours.
15 The postprocessed forecasts of COSMO-DE-EPS in contrast start with lower accuracy but
16 provide meaningful information on longer forecast ranges. The combination of the two sys-
17 tems is performed for probabilities that the expected precipitation amounts exceed a series
18 of predefined thresholds. The resulting probabilistic forecasts are calibrated and outperform
19 both input systems in terms of accuracy for forecast ranges from 1 to 6 hours as shown by
20 verification.

21 The proposed NN-model generalises a previous statistical model based on extended logistic
22 regression, which was restricted to only one threshold of 0.1 mm. The various layers of the
23 NN-model are related to the conventional design elements (e.g. triangular functions and
24 interaction terms) of the previous model for easier insight.

25 1. Introduction

26 Accurate and reliable forecasts of precipitation in the very-short-term range up to 6 h
27 (Wang et al. 2017) in terms of location and time are required in order to issue targeted
28 warnings. Early warnings support to increase the lead time for decision makers in hydrolog-
29 ical and emergency services and can help to diminish possible damages caused by floodings
30 or debris flows. Commonly used bases of these warnings in current operational weather fore-
31 casting are nowcasting systems and numerical weather prediction (NWP). Both approaches
32 can provide valuable warning guidances, however, for different forecast lead times.

33 Methods for nowcasting of precipitation usually rely on remote-sensing observations from
34 a radar network. In a processing step, the obtained radar reflectivities are transformed to
35 estimates of the current rainfall rate. Based on the Lagrangian persistence approach the
36 latest rainfall rates are extrapolated in space and time using a previously determined motion
37 vector field (Germann and Zawadzki 2002). The dynamical evolutions of the convective cells
38 are not considered. Hence, due to the spatial dependency of the lifetime of precipitation
39 fields (Venugopal et al. 1999), such forecasts are able to produce a significant skill as long as
40 the Lagrangian persistence assumption is valid (Zawadzki et al. 1994). The predictability
41 reaches from hours at scales of several hundred kilometers down to minutes when considering
42 developing thunderstorms (Foresti and Seed 2014).

43 In contrast, NWP models explicitly simulate the physical evolution of precipitation fields.
44 Forecast errors result from initial and boundary conditions and from approximating physical
45 equations and their inexact solutions due to finite resolutions in time and space. Especially
46 the simulation of cloud microphysics by sub-grid parameterizations is crucial for precipitation

47 forecasts (Nicolis et al. 2009). In Stephan et al. (2008), it is shown that deficiencies in the
48 simulated rainfall intensities are attributed to shortcomings in the physics parameterizations.

49 In order to reduce the intrinsic uncertainties accompanying NWP forecasts, ensembles were
50 introduced. The ensembles consist of several realizations of a model run, however, differ in
51 initial conditions, boundary conditions and often physical parameterizations (Palmer 2002;
52 Gebhardt et al. 2011). Ensemble forecasting provides users with information on the possible
53 range of weather scenarios to be expected.

54 Nevertheless, multiple runnings of a NWP model can only reduce the random errors and
55 not the aforementioned structural model deficiencies. Therefore, ensemble members are not
56 able to represent the whole spectrum of uncertainties (Scheuerer 2014). Hence, a statistical
57 postprocessing is necessary to reduce systematic biases and an often experienced underdis-
58 persive behavior of ensemble forecasts (Gneiting et al. 2005).

59 Even with a bias corrected and distributional improved NWP prediction, the NWP fore-
60 casts still exhibit various errors in shorter time and small spatial scales. Therefore, a statisti-
61 cal combination with extrapolated nowcasting can improve the skill. The so-called seamless
62 combination aims to create a unique and consistent forecast regardless of location and lead
63 time (Brunet et al. 2015).

64 Vannitsem et al. (2020) point out that the combination of nowcasting and NWP forecasts
65 may basically divided whether they take place in physical or probability space. NIMROD
66 (Golding 1998) as one of the first combination schemes is based on a simple lead-time-
67 dependent weighting function, where the weighting is based on a long-term comparative
68 verification of both initial forecast systems. In INCA (Haiden et al. 2011), the weight for
69 NWP forecasts linearly increases from the beginning until it reaches 1 at a lead time of +4 h.

70 The Short-Term Ensemble Prediction System (STEPS), see Bowler et al. (2006) and Seed
71 et al. (2013), represents a more advanced combination method. Herein, tendencies in the
72 latest observations and the NWP skill are quantified in real time and used to adjust weights
73 combining the nowcast extrapolation and the NWP forecast, in dependence on lead time
74 and spatial scale. Additionally, a forecast ensemble is generated due to the replacement
75 of non-predictable spatial scales with correlated stochastic noise. Due to the emerging
76 of nowcasting ensembles, efforts were made to not only use the forecast skill as objective
77 combination metric but also the ensemble spread. Recently, Nerini et al. (2019) utilize an
78 ensemble Kalman filter to iteratively combine NWP forecasts with radar-based nowcasting
79 extrapolations.

80 In reference to combination methods in probability space, the blending scheme of Kober
81 et al. (2012) weights exceedance probabilities derived from output of the NWP model
82 COSMO-DE-EPS with smoothed neighborhood probabilities computed from the determin-
83 istic nowcasting algorithm Rad-TRAM. The weights are based on a long-term verification.
84 Combination of multi-model ensembles are carried out in Johnson and Wang (2012) and
85 Bouttier and Marchal (2020).

86 In a previous study, Schaumann et al. (2020) propose the LTI-model as a modified logistic
87 regression model for precipitation rates higher than 0.1 mm h^{-1} . Additionally to the logistic
88 regression, triangular functions and interaction terms are introduced to take the possible
89 differences in the individual initial probabilistic forecasts into account and, furthermore, to
90 increase the flexibility to compensate a possible underestimation and overestimation. The
91 logistic regression model is a common tool in the area of probabilistic weather forecasting
92 and has been used for the calibration of forecasts in ,e.g., Hamill et al. (2008).

93 The present study aims to generalize the LTI-model with the help of a neural network
94 (NN). The network to be developed should not only satisfy the demands set for the LTI-
95 model but should also provide consistent threshold exceedance probabilities, where con-
96 sistency is understood such that the probabilities of exceeding a threshold monotonously
97 decrease with increasing thresholds. This is not ensured when models are trained for each
98 threshold independently. Further, the network should be able to represent forecast uncer-
99 tainty with increasing lead time. For other extensions to the logistic regression model in
100 order to ensure consistency, see Wilks (2009); Ben Bouallègue (2013).

101 As in Schaumann et al. (2020), the training data set is based on forecasts of RadVOR
102 (Winterrath et al. 2012) and Ensemble-MOS (Hess 2020). RadVOR is a nowcasting system
103 that provides deterministic extrapolation forecasts of radar-based rainfall estimates. Radar
104 observations are obtained by the operational German radar network operated by Deutscher
105 Wetterdienst (DWD). Exceedance probabilities from the deterministic extrapolation fore-
106 casts are derived by using the neighborhood approach described by Theis et al. (2005).
107 Ensemble-MOS statistically postprocesses output from the ensemble of the convection-
108 permitting COSMO-DE model, which was upgraded to COSMO-D2 on 15 May 2018. It
109 provides probabilistic precipitation forecasts for various thresholds using logistic regression.

110 The present study is organized as follows. In Section 2, a brief overview of the utilized
111 training data set is given. Section 3 provides a brief summary of the LTI-model. Afterwards,
112 the development of a NN architecture for the model generalization regarding the simulta-
113 neous consideration of multiple thresholds is described in detail. Since NNs react sensitive
114 on the chosen hyper-parameters, a hyper-parameter optimization approach is provided in
115 the appendix, see Section A1. Results are stated in Section 4. Herein, sensitivities in the

116 choice of hyper-parameters are investigated and a combination example is given. Finally, in
117 Section 5 conclusions are drawn and an outlook is given for possible future developments.

118 **2. Data**

119 As in Schaumann et al. (2020), we use forecasts of the DWD systems Ensemble-MOS and
120 RadVOR as data sources. However, we now extend the considered training and forecast
121 period by three months and have thus, altogether, 6 months of precipitation data from
122 April to September 2016. Furthermore, in the previous study, we considered data for only
123 one threshold (0.1mm h^{-1}) whereas now we consider 9 precipitation thresholds $t_1 = 0.1$,
124 $t_2 = 0.2$, $t_3 = 0.3$, $t_4 = 0.5$, $t_5 = 0.7$, $t_6 = 1$, $t_7 = 2$, $t_8 = 3$ and $t_9 = 5$, with mm h^{-1} as unit
125 of measurement.

126 The data used from both sources for the results derived in the present paper span across
127 Germany and parts of neighboring countries. In this section, we briefly recall the main
128 properties of the forecast systems Ensemble-MOS and RadVOR as well as the calibrated
129 rainfall estimates used as ground truth. For further details regarding this data set, we refer
130 to our previous study.

131 *a. Ensemble-MOS*

132 The postprocessing system Ensemble-MOS of DWD is a model output statistics (MOS)
133 system with the capability to statistically process ensemble forecasts resulting in hourly
134 outputs of probabilistic forecasts. These forecasts are available on a regular grid of 20
135 20 km^2 , for lead times up to +21 h, and for various weather elements to support weather
136 warnings (Hess 2020). From the latter output variables, we consider the exceedance of
137 nine precipitation thresholds. The Ensemble-MOS forecasts used in the present paper are

138 based on ensemble forecasts of DWD’s 2016 operational convection-permitting NWP model
139 COSMO-DE-EPS (Theis et al. 2012). In particular, the training of the Ensemble-MOS relies
140 on COSMO-DE-EPS forecasts from 2011 to 2015.

141 *b. RadVOR*

142 Additionally to the Ensemble-MOS forecasts, we use data from the nowcasting method
143 RadVOR (Weigl and Winterrath 2010; Winterrath et al. 2012). RadVOR provides every
144 5 minutes deterministic precipitation forecasts for lead times up to +2h on a regular grid
145 of size $1 \times 1 \text{ km}^2$. These very-short term forecasts consist of two components. In a first
146 step, estimates of the current rainfall rate are derived from radar reflectivities. Thereafter,
147 these rainfall rates are then advected in 5 minute increments according to a previously esti-
148 mated motion vector field. Note that the edges of the respective forecast domain are shifted
149 accordingly.

150 For the combination with Ensemble-MOS using the model proposed in the present paper,
151 the RadVOR forecasts are interpolated to the same grid and matching time intervals. For
152 this, the 5-minute precipitation amounts are added up for a one-hour sum. Next, we consider
153 the grid points with precipitation rates larger than the given threshold on the $1 \times 1 \text{ km}^2$
154 grid. Finally, to estimate the probability of threshold exceedance, we compute the average
155 exceedance on the $1 \times 1 \text{ km}^2$ grid for a neighbourhood of each grid point on the $20 \times 20 \text{ km}^2$
156 grid.

157 *c. Calibrated hourly rainfall estimates*

158 As a ground truth for the training and validation of our models, we use calibrated rainfall
159 estimates on a $1 \times 1 \text{ km}^2$ grid based on reflectivity measurements of DWD’s operational radar

160 network. These rainfall estimates are adjusted by about 1,300 rain gauge measurements to
161 reduce the error induced by the uncertainty of the relationship between radar reflectivities
162 and precipitation amounts (Winterrath et al. 2012). Note that, in comparison to the previous
163 study, the filter algorithm proposed by Winterrath and Rosenow (2007) for removing pixel
164 artifacts has not been applied.

165 **3. Neural Network Architectures**

166 In Schaumann et al. (2020) we used a modified logistic regression model for the combi-
167 nation of two different probabilistic forecasts. This model is referred to as LTI-model in
168 the following. Here, L stands for "logistic" while T and I refer to "triangular functions"
169 and "interaction terms", respectively, which are the modifications of the standard logistic
170 regression, where the latter one is called the L-model.

171 To begin with, we briefly summarize the basic idea of the LTI-model and, later on in
172 Section 3b, we propose further modifications of it for the simultaneous consideration of
173 multiple thresholds.

174 *a. Modified logistic regression model*

175 The introduction of the LTI-model aimed to develop a model for the seamless calibrated
176 combination of the afore described forecast systems Ensemble-MOS and RadVOR, see Sec-
177 tion 2. Here, statistical calibration refers to an ideal reliability diagram, which is a desirable
178 property of probabilistic forecasts (Murphy and Winkler 1977, 1987). The combined fore-
179 cast should outperform the individual initial forecasts for all relevant lead times in accuracy
180 and reliability. For short lead times the extrapolated nowcasting of RadVOR outperforms
181 the skill of Ensemble-MOS. However, with increasing lead time, forecast scores for RadVOR

182 drop rapidly in comparison to those for Ensemble-MOS, see Figure 1. Note that the green
 183 plots in Figure 1 illustrate the notable improvements achieved by the NN-model introduced
 184 in Section 3b below.

185 The LTI-model is based on the standard logistic regression model $f_L : [0;1]^n \rightarrow [0;1]$ for
 186 some $n > 1$, where n denotes the number of probabilistic input forecasts to be combined.
 187 From a mathematical point of view, the L-model estimates the conditional probability distri-
 188 bution of a dichotomous random variable $Y : \Omega \rightarrow \{0,1\}$, i.e. the occurrence of precipitation
 189 above a certain fixed threshold, conditioned on given realizations x_i of a family of random
 190 variables $X_i : \Omega \rightarrow [0;1]$ for $i \in \{1, \dots, n\}$, i.e. the probabilistic input forecast models.
 191 The random variables $Y; X_1; \dots; X_n$ are defined on some probability space $(\Omega; \mathcal{F}; \mathbb{P})$, where
 192 Ω contains the vectors $(y; x_1; \dots; x_n)$ consisting of all possible forecasts $x_1; \dots; x_n$ of the n
 193 input forecast models $X_1; \dots; X_n$, and the precipitation occurrence indicator y . Thus,

$$f_L(x_1; \dots; x_n) = \mathbb{P}(Y = 1 \mid X_1 = x_1; \dots; X_n = x_n) \quad \text{for all } x_1; \dots; x_n \in [0;1].$$

194 Note that in Schaumann et al. (2020) the special case $n = 2$ has been considered, where
 195 the probabilities $x_1; x_2$ originate from RadVOR and Ensemble-MOS, respectively. Then, the
 196 L-model combines these two input forecasts and provides calibrated forecast probabilities.

197 1) Triangular functions

198 For the case $n = 2$ the L-model has three weights $w_0; w_1; w_2 \in \mathbb{R}$ and is given by

$$f_L(X_1; X_2) = \sigma(w_0 + w_1 X_1 + w_2 X_2); \tag{1}$$

199 where $\sigma(x) = \frac{e^x}{e^x + 1}$ is the logistic function. Note that there are individual weights for each
 200 forecast time step. Due to the small number of weights, the L-model is rather limited in
 201 how it combines the input forecast models X_1 and X_2 . The weights merely allow for enough

202 flexibility to weight each forecast based on its overall forecast quality. However, the forecast
 203 quality of X_i might vary across the range of possible predictions within the interval $[0;1]$.
 204 This variation in forecast quality is expressed by the reliability diagram of X_i (see Figure 2
 205 for examples) indicating for which values $x_i \in [0;1]$ the input forecast model X_i tends to
 206 over- or underestimate the occurrence of the event that $Y = 1$. Therefore, we proposed to
 207 choose the weights for each model $X_1; X_2$ in dependence of the forecast values x_1 and x_2 .
 208 For this, a family of $m + 1$ triangular functions $f_j : [0;1] \rightarrow [0;1]$ has been defined for some
 209 integer $m > 0$ and all $j = 0; \dots; m$ with

$$f_j(x) = \max\left\{0; \frac{j - mx}{m}\right\} \quad \text{for all } x \in [0;1] \quad (2)$$

and $\sum_{j=0}^m f_j(x) = 1$ for $x \in [0;1]$. Thereby, for $i = 1; 2$, the input forecast model X_i can be
 encoded as a random vector $(X_i) = (f_0(X_i); \dots; f_m(X_i)) \in [0;1]^{m+1}$. Thus, at most two
 consecutive triangular functions $f_j(X_i)$ and $f_{j+1}(X_i)$ are non-zero, which is the case when
 the value of X_i falls between $\frac{j}{m}$ and $\frac{j+1}{m}$. Now, instead of the forecast models X_1 and X_2 ,
 we pass the random vectors (X_1) and (X_2) to the L-model and call that the LT-model.
 For some weights $w_{ij} \in \mathbb{R}$, the latter model is defined as

$$f_{\text{LT}}(X_1; X_2) = f_{\text{L}}((X_1); (X_2)) = \sum_{i=1}^2 \sum_{j=0}^m w_{ij} f_j(X_i)$$

210 If we now consider a reliability diagram with $m + 1$ bins, then each weight for a triangular
 211 function corresponds to one bin and, therefore, allows the model to produce a calibrated
 212 forecast. Note that the LT-model does not use a weight w_0 like in Eq. (1), sometimes called
 213 “bias“ or “intercept“, because w_0 is redundant as the triangular functions sum up to 1 for
 214 any $x \in [0;1]$.

215 2) Interaction terms

216 The weights in the L- and LT-models are chosen for either one of the two input fore-
 217 cast models $X_1; X_2$, irrespective of the other forecast model. However, it might be sensible
 218 to choose different weights, depending on whether both forecasts agree or disagree on the
 219 probability of occurrence of the event $Y = 1$. For this, we consider four additional pre-
 220 dictors ${}_1(X_1; X_2); \dots; {}_4(X_1; X_2)$, called interaction terms, where ${}_1(X_1; X_2) = \mathbb{P}(\overline{X_1} \overline{X_2})$,
 221 ${}_2(X_1; X_2) = \mathbb{P}(\overline{(1 - X_1)} X_2)$, ${}_3(X_1; X_2) = \mathbb{P}(\overline{X_1} (1 - X_2))$, ${}_4(X_1; X_2) = \mathbb{P}(\overline{(1 - X_1)(1 - X_2)})$.

222 Like in the previous section, where we have considered the vectors $(X_i) =$
 223 $({}_0(X_i); \dots; {}_m(X_i))$ for $i = 1; 2$, we now apply triangular functions to the interaction
 224 terms and pass the random vectors $({}_i(X_1; X_2)) = ({}_0({}_i(X_1; X_2)); \dots; {}_m({}_i(X_1; X_2)))$ for
 225 $i = 1; 2; 3; 4$ to the model, i.e.,

$$\begin{aligned} f_{\text{LTI}}(X_1; X_2) &= f_{\text{L}}((X_1); (X_2); ({}_1(X_1; X_2)); \dots; ({}_4(X_1; X_2))) \\ &= \sum_{i=1}^4 \sum_{j=0}^m w_{ij} {}_j(X_i) + \sum_{i=1}^4 \sum_{j=0}^m w'_{ij} {}_j({}_i(X_1; X_2)) ; \end{aligned}$$

226 where $w_{ij}; w'_{ij} \in \mathbb{R}$ are some weights.

227 3) Model training

228 For training the LTI-model, a rolling-origin with reoptimization scheme (Armstrong and
 229 Grohman 1972) is used in order to simulate the operational conditions. For this, the available
 230 data is not split up in separate training and test datasets, but it is split by a point in time
 231 t , which represents the “present“, into “past data“ and “future data“. In each step of
 232 the rolling-origin scheme the model is trained on “past data“ and then validated based on
 233 predictions made for time steps ahead of t , on which the model has not been trained yet. At

234 the end of each step, τ is moved forward in time by one time step. This process is repeated
235 until τ traversed the entire dataset.

236 *b. Generalization of the LTI-model for multiple thresholds, using neural networks*

237 1) Overview

238 In the present paper, we propose a number of modifications of the LTI-model, which allow
239 for the combination of forecasts for several thresholds by one common model. Note that the
240 LTI-model, being a modified logistic regression model, can be seen as a specific type of a
241 NN. Therefore, more general variants of NNs are a natural choice to make further extensions
242 of the LTI-model.

243 In NNs two types of parameters are distinguished: hyper-parameters and trainable pa-
244 rameters. Hyper-parameters determine the architecture of the NN-model (e.g. the numbers
245 and types of layers and neurons). They have to be determined before fitting the trainable
246 parameters to a data set. The trainable parameters are the weights within each layer. The
247 fitness of a specific architecture as defined by the hyper-parameters is highly dependent on
248 the problem to be solved. While there are some guidelines on how to design a NN, it is impos-
249 sible to tell how the choice of a hyper-parameter affects the performance before fitting and
250 validating the model (Bergstra et al. 2011). For the optimization of the hyper-parameters
251 of our NN-model, we propose an algorithm in Section A1.

252 In the following sections we give a brief overview of the components of the proposed NN-
253 model. Each component is either a generalization of a part of the LTI-model or a newly
254 added modeling component. For each of them we introduce a number of hyper-parameters,
255 which define the architecture of the NN-model and which have to be determined before
256 training of the model can begin. In total there are 10 hyper-parameters. Each hyper-

257 parameter has a set of possible configurations, which we denote by $H_i = f c_i^1; \dots; c_i^{m_i} g$ for
 258 $i = 1; \dots; 10$. Thus, the architecture of the NN-models considered in this paper is defined
 259 by a vector $h \in H_1 \times \dots \times H_{10}$, which contains exactly one possible configuration for each
 260 hyper-parameter. For each $h \in H_1 \times \dots \times H_{10}$, the NN-model consists of the following
 261 layers: (i) Zero to five convolutional layers, (ii) Zero or one dense layer for interaction
 262 terms (iii) Zero or one layer consisting of triangular functions, (iv) one softmax layer. For
 263 a schematic representation of the model architecture, see Figure 3. With the exception of
 264 convolutional layers, each layer in this NN architecture corresponds to a component of the
 265 LTI-model, whereby the dense layer and softmax layer are more general than their LTI-
 266 counterparts. For an introduction to deep learning and explanations regarding topics like
 267 “activation function“ or “convolutional layer“, see Chollet (2017).

268 2) Convolutional layers for model input

269 With increasing lead time, it becomes more difficult to make accurate predictions due
 270 to increased forecast uncertainty. In particular, this might cause forecast models to be
 271 imprecise in the prediction of the location of precipitation events. For probabilistic forecasts,
 272 this imprecision may manifest itself in two ways: (i) The precipitation events are predicted
 273 at wrong locations. Note that this is especially relevant for the RadVOR probabilities,
 274 which are based on a simple extrapolation and therefore do not take increased uncertainty
 275 for longer lead times into account. This results in equally sharp predictions for all lead
 276 times. (ii) The probability mass spreads out spatially. In other words, the forecast model
 277 predicts lower precipitation probabilities for a larger area to take the increased uncertainty
 278 into account.

279 In both cases, it is advantageous for a combination model to be aware of the predictions
280 made at adjacent locations in order to combine two forecasts at a given location. For this,
281 we consider convolutional layers (Zhang et al. 1990), which are commonly used for analyzing
282 image data or data arranged on a regular grid. A NN with convolutional layers takes the
283 information from a neighborhood of adjacent grid points into account, whereas the LTI-
284 model combines forecasts point-wise, i.e., the model output for each grid cell depends only
285 on the individual input forecasts at this location. The size of this neighborhood is determined
286 by the sizes of the convolutional kernels, which map the neighborhood data onto a vector of
287 a specified length.

288 In the present paper, all kernels are square-shaped and therefore the considered sizes refer
289 to both height and width, e.g., a kernel size of 3 refers to a neighborhood of 3 × 3 grid cells.
290 Both the size of the kernel and the length of its output are hyper-parameters of the layer.
291 The input of a convolutional layer is a tensor $I \in \mathbb{R}^{b_x \times b_y \times b_z}$, where b_x and b_y define the size
292 of the forecast grid in x - and y -direction, respectively. Furthermore, b_z depicts the number
293 of probabilistic predictions of Ensemble-MOS and RadVOR for 9 thresholds each. Thus,
294 $b_z = 18$.

295 For technical reasons, the NN requires input data given on a rectangular grid. In cases
296 where some of the grid cells are undefined (i.e. where no forecast is available), the input data
297 is restricted to the largest rectangle, which contains only valid values, i.e., is NaN-free. Since
298 the forecasts of RadVOR are based on radar measurements that are extrapolated according
299 to a motion vector field, the edges of the respective forecast domain shift in time. Thus,
300 the area containing data of both RadVOR and Ensemble-MOS depends on the magnitude
301 of the motion vector field and may decrease with lead time. Another limitation is the total
302 convolutional size, which has to fit inside the well-defined rectangle.

303 The total convolutional size is given by the formula: $c_1(c_2 - 1) + 1$ with $c_1 \geq H_1$ (number of
304 convolutional layers), $c_2 \geq H_2$ (kernel size of each layer). Due to this, some combinations of
305 H_1 and H_2 result in an oversized convolution and, therefore, only combinations with a total
306 convolutional size less than or equal to 13 are used. The latter convolution corresponds to
307 a neighborhood of 260 km.

308 The following hyper-parameters and configurations are considered: number of convolu-
309 tional layers with $H_1 = f0;1;2;3;4;5g$, kernel size with $H_2 = f1;2;3;5;6;7;9;11;13g$,
310 length of the output vector, i.e., the number of convolution matrices used by the
311 layer (Chollet 2017), with $H_3 = f1;2;4;6;8;10;12;14;16g$, activation functions with
312 $H_4 = ff_{\text{elu}};f_{\text{exp}};f_{\text{lin}};f_{\text{sigmoid}};f_{\text{relu}};f_{\text{tanh}}g$ and L₂-regularization strength with $H_5 =$
313 $f0;10^{-7};5 \quad 10^{-7};10^{-6};5 \quad 10^{-6};10^{-5}g$.

314 3) Dense layer for interaction terms

315 For the LTI-model, the interaction terms were chosen by hand prior to model fitting. In the
316 framework of NNs, the functionality of the interaction terms can be achieved with a densely
317 connected layer of neurons, see e.g. Chollet (2017). In comparison to the LTI-model, a dense
318 layer has the advantage that the shapes of the interaction terms are mostly determined by
319 trainable parameters. The hyper-parameters for this layer are: number of neurons with
320 $H_6 = f0;2;4;6;8;10;12g$, activation functions with $H_7 = ff_{\text{elu}};f_{\text{exp}};f_{\text{lin}};f_{\text{sigmoid}};f_{\text{relu}}g$ and
321 L₂-regularization strength with $H_8 = f0;10^{-7};5 \quad 10^{-7};10^{-6};5 \quad 10^{-6};10^{-5}g$.

322 4) Layer for triangular functions

323 The role of this layer is the integration of triangular functions into the NN-model.
324 For a definition of and the rationale behind triangular functions, see Section 3a. The

325 only hyper-parameter for this layer is the number of triangular functions with $H_9 =$
 326 $\{0; 2; 4; 6; 8; 10; 12; 14\}$. Thus, for each $c \in H_9 \setminus \{0\}$, this layer applies the triangular func-
 327 tions $f_0; \dots; f_c$ to the output of each neuron of the previous dense layer. In our case, for
 328 each $c' \in H_6$, the previous layer returns a vector of scalars $(x_1; \dots; x_{c'})$. Hence the triangular
 329 functions layer has the output $(f_0(x_1); \dots; f_c(x_1); \dots; f_0(x_{c'}); \dots; f_c(x_{c'}))$ of size $(c+1)c'$. To
 330 our knowledge, such functions are not used in conventional NNs (a similar concept are radial
 331 basis function networks (Park and Sandberg 1991)), but they can be manually constructed
 332 with the help of Keras backend functions (Keras 2020).

333 5) Softmax layer to ensure consistent predictions for all thresholds

334 To obtain exceedance probabilities for all considered thresholds $t_1; \dots; t_m$ with $0 < t_1 <$
 335 $\dots < t_m$ by means of LTI-models, a separate LTI-model would have to be trained for
 336 each threshold. However, this does not guarantee that the probabilities are decreasing
 337 monotonously for increasing thresholds, since the separate LTI-models have no knowledge
 338 about each other. Hence, an extended model is needed, of which the output is a vector
 339 of monotonously decreasing probabilities for the exceedance of the considered thresholds
 340 $t_1; \dots; t_m$. Additionally, the data available for one threshold might be useful for the combi-
 341 nation of other thresholds, too, since each probability is a point on the discrete cumulative
 342 distribution function of the same event and, therefore, they are interlinked.

343 To ensure that the components of the vector of combined forecasts are decreasing
 344 monotonously, we train the neural network on a multi-label classification problem with
 345 a Softmax layer (Bridle 1990). This can be seen as a generalization of the logistic re-
 346 gression model, which has been utilized in the LTI-model. While the logistic regression
 347 model estimates the conditional probability distribution of a dichotomous random variable

348 $Y : \Omega \rightarrow [0; 1]g$, the softmax layer allows for estimating the conditional probability distri-
 349 bution of a random variable $Y' : \Omega' \rightarrow [0; \dots; m]g$ where $m \geq 1$. In our case, Y' models
 350 the exceedance of the considered precipitation thresholds at a given location. For this, let
 351 $T : \Omega' \rightarrow [0; 1)$ be the precipitation amount and let $C_i = \{T \geq [t_i; t_{i+1})\}g$ denote the event
 352 that T takes values between the thresholds t_i and t_{i+1} , for $i = 0; \dots; m$. For this, we formally
 353 introduce two further thresholds t_0 and t_{m+1} , where $t_0 = 0$ and $t_{m+1} = 1$. We then put
 354 $Y' = i$ if $T \in [t_i; t_{i+1})$, i.e., Y' indicates which of the events $C_0; \dots; C_m$ is occurring.

355 For the family of pairwise disjoint events $C = \{C_0; \dots; C_m\}g$, the NN learns to predict
 356 a conditional discrete probability distribution $P_{C;I} = \{P(C_0 | I); \dots; P(C_m | I)\}g$, where
 357 $P(C_i | I)$ is the conditional probability for the occurrence of event C_i given the model input
 358 I . Then, for each $j \in [1; \dots; m]g$, the conditional probability of the event that $T \geq t_j$ can be
 359 computed by $P(T \geq t_j | I) = 1 - \prod_{i=0}^{j-1} P(C_i | I)$. Clearly, from this it follows by definition
 360 that $P(T \geq t_j | I) \geq P(T \geq t_k | I)$ if $t_j < t_k$, and therefore the predictions of the NN are
 361 consistent for all thresholds.

362 6) Optimizer for trainable parameters

363 Another difference between the LTI-model introduced in Schaumann et al. (2020) and the
 364 NN-model considered in the present paper is the choice of the optimizer. Note that the
 365 optimizer controls how the trainable parameters change during the model training. This
 366 has a large influence on the model fitness. For a more detailed introduction to the operating
 367 principle of optimizers, we refer to Chollet (2017).

368 In our previous paper, a stochastic gradient descent with a constant learning rate has been
 369 used as optimizer for the LTI-model. While this is sufficient for the (relatively small) number
 370 of parameters of the LTI-model, the NN-model considered in the present paper requires a

371 more sophisticated optimizer, due to weaknesses of the classical stochastic gradient descent.
372 Depending on the network architecture and the training data set, gradients for some weights
373 might “vanish“ at some point in training, see Glorot et al. (2011). This means that parts of
374 the NN might receive only small updates or a sparse number of updates leading to stagnation
375 of the training process. More recent optimizers address this problem by various means, e.g.,
376 gradient descent with momentum (Sutskever et al. 2013) or adaptive learning rates. This
377 ensues, first, to scale up small gradients back to a reasonable size or, second, to compensate
378 for a sparse number of updates of a weight. In the present paper, five different optimizers are
379 investigated. All of them are based on adaptive learning rates, which leads to a tenth hyper-
380 parameter with $H_{10} = fAdam; Adagrad; Adadelata; Adamax; Nadamg$. For more details on
381 how each optimizer works, see Kingma and Ba (2015), Duchi et al. (2011), Zeiler (2012),
382 and Dozat (2016).

383 *c. Training and Validation*

384 In the previous section, several modifications of the LTI-model have been discussed. Each
385 of them introduces one or more hyper-parameters, needing to be determined before the NN
386 can be trained. For this, a hyper-parameter optimization algorithm is employed, which is
387 explained in Section A1 in more detail.

388 To train the NN, the available data is split into several parts. One for each step of the
389 training process: (i) Training of different model architectures for the hyper-parameter search,
390 using data of the period from 2016-04-01 to 2016-05-31, (ii) fitness evaluation of model
391 architectures for the hyper-parameter search, using data from 2016-06-01 to 2016-06-30,
392 (iii) pick best architecture based on evaluation results, (iv) training of the best architecture

393 (without validation), using data from 2016-04-01 to 2016-06-31, and (v) rolling origin update
394 and validation with best architecture, using data from 2016-07-01 to 2016-09-30.

395 This kind of training process ensures that the choice of the best performing architecture
396 and its validation are based on two different time intervals, since this would lead to a skewed
397 validation result otherwise.

398 4. Results

399 *a. Influence of hyper-parameters*

400 1) Optimizer

401 For the hyper-parameter optimization, about 18000 model architectures have been eval-
402 uated for each considered lead time. The choice of the optimizer $c \in H_{10}$ has, by far, the
403 largest impact on model fitness (see Section A1a). The distribution of model fitness for each
404 optimizer is depicted in Figure 4. Since Adam, Adamax and Nadam outperform Adagrad
405 and Adadelata for almost all model configurations, we will focus on results only with regard
406 to Adam, Adamax and Nadam in the following discussion.

407 2) Convolutional layers

408 Furthermore, the model fitness is highly affected by the number of convolutional layers
409 (selected from the set H_1) and their kernel size (selected from H_2), see Figure 5, where the
410 results of the hyper-parameter optimization are visualized. While the difference between
411 individual configurations is less pronounced for lead times +1 h, larger total convolutional
412 sizes perform better for longer lead times. Thus, in situations of increased forecast uncer-
413 tainty (e.g. for longer lead times), an improved forecast skill is achieved when considering

414 more adjacent grid cells. This is in agreement with the ideas of Theis et al. (2005) and
415 Schwartz and Sobash (2017).

416 The activation functions f_{elu} ; f_{linear} ; f_{relu} and f_{tanh} seem to perform similarly well when
417 compared to each other. In contrast, the functions $f_{\text{exponential}}$ and f_{sigmoid} exhibit a much
418 worse behavior, in particular for model architectures with many convolutional layers. As
419 an exception, f_{sigmoid} does not show this behavior for the lead time +6 h and even performs
420 best out of all considered activation functions.

421 Models with larger lengths of output vectors (selected from H_3) tend to perform better,
422 however, the difference is clearly pronounced only up to a vector length of 4.

423 It should be noted that the number of convolutional layers (selected from H_1) and their
424 kernel size (selected from H_2) affect the output size of the NN in two different ways: (i)
425 As explained in Section 3b, input data passed to the NN needs to be rectangular-shaped
426 and defined on each grid cell in order to enable the NN to learn and to make a prediction.
427 Note that the passed data domain underlies a large variability, due to unregular boundaries
428 of and occasionally missing data within the input forecasts. The total convolutional size
429 determines the possible minimum edge length of the data domain. (ii) Furthermore, the
430 total convolutional size determines the size of the input area which is mapped to an output
431 value. For example, for a total convolutional size of 5, a model input I of size 15 30 18
432 is mapped to a model output of size 11 26 9.

433 Due to these effects, the amount of available data used for training and validation depends
434 on the total convolutional size. This might be an explanation for the decreased fitness of the
435 four configurations (1;13);(4;4);(3;5);(2;7) of H_1 H_2 with the largest total convolutional
436 size of 13, since they have much less output to be trained and validated on, see Figure 5.

437 For a total convolutional size of 9 (used for lead times +1 h, +2 h, +4 h and +5 h, see
438 Table 1), the passed data domain I consists, on average, of about 292 grid cells in the
439 considered time period (July, August, September 2016). This results in 636126 data points
440 in total. In comparison to this, for a total convolutional size of 11 (used for lead times +3 h
441 and +6 h, see Table 1), the passed data domain consists, on average, of about 160 grid cells
442 with 344431 data points in total.

443 3) Triangular functions and interaction terms

444 In Figure 6 the effect of triangular functions and the dense layer (interaction terms) on
445 the model fitness is visualized. In general, one might expect that more neurons in the
446 dense layer perform better than less. Thus, at first glance, it seems to be counter-intuitive
447 that 2 neurons perform worse than no neurons at all. A likely explanation of this is that
448 few neurons act as a bottleneck that restricts the amount of information the NN can pass
449 through the dense layer, whereas for zero neurons the layer and, therefore, the bottleneck
450 is removed. Regarding triangular functions, typically 3 to 9 of them perform best for all
451 lead times. However, some of these configurations can perform worse for specific lead times,
452 e.g., 11 triangular functions for +4 h, and 5 triangular functions for +5 h, see Figure 6. The
453 results visualized in Figure 6 show that model fitness for lead time +1 h behaves differently
454 in comparison to the model fitnesses for longer lead times. See also Figure 5, where this
455 effect can be observed, too.

456 4) Remaining hyper-parameters

457 The remaining hyper-parameters are the activation functions (selected from H_7) for the
458 dense layer and the L_2 regularization strengths for the convolutional layers (selected from

459 H_5) and the dense layer (selected from H_8). However, these parameters do not seem to affect
460 the model fitness in any significant way.

461 *b. Validation of chosen architectures for different lead times*

462 The model fitness has been evaluated on data for the months July, August and September
463 2016. Within this period of time and the passed domain I (depending the total convolu-
464 tional size), the considered precipitation thresholds were exceeded as described below by the
465 numbers of upcrossings (and corresponding relative frequencies in parenthesis).

466 (i) For a total convolutional size of 9 we have: 39303 (6.18%) for 0.1 mm, 31284 (4.92%)
467 for 0.2 mm, 26402 (4.15%) for 0.3 mm, 20300 (3.19%) for 0.5 mm, 16459 (2.59%) for
468 0.7 mm, 12677 (1.99%) for 1 mm, 6264 (0.98%) for 2 mm, 3494(0.55%) for 3 mm, and
469 1375 (0.22%) for 5 mm.

470 (ii) For a total convolutional size of 11 we have: 20420 (5.93%) for 0.1 mm, 16195 (4.70%)
471 for 0.2 mm, 13561 (3.94%) for 0.3 mm, 10330 (3.00%) for 0.5 mm, 8358 (2.43%)
472 for 0.7 mm, 6386 (1.85%) for 1 mm, 3046 (0.88%)for 2mm, 1652 (0.48%) for 3 mm,
473 677 (0.20%) for 5 mm.

474 For each lead time, the parameter configurations are shown in Table 1, which were chosen
475 by the hyper-parameter optimization algorithm explained in Section A1. For each chosen
476 architecture validation scores and reliability diagrams are shown in Figures 1 and 2. Note
477 that for some hyper-parameters several configurations perform equally well, which leads to
478 random fluctuations between the choices for different lead times.

479 It can be seen that the NN-model makes less biased and more calibrated predictions with
480 a higher Brier skill score in comparison to both initial forecasts provided by Ensemble-MOS

481 and RadVOR, for all lead times and thresholds. Although the RadVOR forecasts are only
482 provided up to +2 h, the combined forecasts have better scores up to +6 h. A similar effect
483 has been observed in our previous paper regarding the scores of combined forecasts obtained
484 by the LTI-model. For more details on the used validation scores, see Wilks (2006).

485 Similar to the previously considered LTI-model, the NN-model of the present paper has
486 improved reliability diagrams in comparison to both initial forecasts, see Figure 2. As
487 expected, the reliability of forecasting models decreases with increasing lead times and in-
488 creasing thresholds. In order to keep the reliability diagrams calibrated, the NN learns to
489 lower its predictions accordingly to not overestimate the occurrence of precipitation, which
490 leads to shorter curves in the reliability diagram for longer lead times and higher thresholds
491 in comparison to both initial input models.

492 To test if it is actually necessary to run the hyper-parameter optimization algorithm for
493 each lead time, we trained the chosen architecture for +1 h on the other five lead times, too.
494 A visual comparison between the validation scores given in Figure 1 and the results for the
495 +1 h-architecture showed only slight fitness differences. However, the reliability diagrams
496 seem to be less calibrated, see Figure 2, We therefore decided to use a separate network
497 architecture for each lead time.

498 *c. Combination example*

499 In Figure 7, forecasting results obtained by the combination of input data from Ensemble-
500 MOS and RadVOR are shown, as an example, for the hour 6-7 am of 2016-07-21 and for
501 three lead times (+1 h, +2 h, +3 h). Due to a larger total convolutional size of the NN-
502 model for +3h, the size of the output is smaller. For shorter lead times, the forecast of the

503 NN-model closely resembles the forecast of RadVOR, while for increasing lead times, the
504 predictions become more smooth and more dependent on the Ensemble-MOS prediction.

505 **5. Conclusions**

506 *a. Summary of results*

507 In this paper we presented a hyper-parameter optimization algorithm and NN architectures
508 for the combination of two probabilistic forecasts, where we consider several precipitation
509 thresholds simultaneously. The architectures chosen by the hyper-parameter optimization
510 algorithm show improvements for all considered validation scores across all thresholds, and
511 calibrate the resulting probabilities. Like the previously developed LTI-model, the NN-
512 model considered in the present paper improves forecast scores also for lead times longer
513 than +2h, although RadVOR forecasts were only provided up to +2h.

514 When comparing forecast scores of the LTI-model and the NN-model, it seems that the
515 reliability diagrams of the LTI-model are smoother than those of the NN-model. This might
516 be due to the fact that in the training process of the NN the training data must be cropped
517 down to a smaller rectangular area with valid forecast values, which discards data points for
518 the NN-model that are still available to the LTI-model. Also, the NN must predict consistent
519 exceedance probabilities for several thresholds, which is an additional constraint to satisfy.
520 It might be necessary for the NN-model to sacrifice some reliability for one threshold, in
521 order to make better predictions for all considered thresholds on average.

522 *b. Outlook & possible next steps*

523 According to the results of the hyper-parameter search performed in the present paper,
524 some hyper-parameters seem to be much more important than others. Thus, it might

525 be possible to further improve the architecture by adapting the search space. Since the
526 optimizer and the number and size of convolutional layers have the largest influence on model
527 fitness, additional optimizers and convolutional layer combinations should be investigated.
528 Thus, in a forthcoming paper, we will investigate the numerical stability of the hyper-
529 parameter optimization algorithm and how the chosen architectures in Table 1 compare
530 ,e.g., to the second best architecture for each lead time.

531 Due to the restriction that the input of the NN must be rectangular and free of missing
532 values, it should be considered to generate valid values by interpolation at grid points with
533 missing values, or to pass an mask to the NN as additional input in order to specify, which
534 values are valid. This would allow training without cropping of the data and also increase
535 the area for which predictions can be made.

536 Furthermore, it would be interesting to investigate how additional information might affect
537 the quality of combination, e.g., by increasing the resolution of the grid, by passing ensemble
538 members directly to the NN without aggregation to probabilities, by adding an orography
539 map to the input, or by including additional meteorological indicators. Given that a new
540 dataset contains enough precipitation events for higher thresholds, the list of considered
541 thresholds could be expanded.

543 **A1. Hyper-parameter optimization**

544 To choose hyper-parameters by means of a systematic approach, various optimization algo-
545 rithms have been developed, attempting to find correlations between the hyper-parameters
546 of a model and its fitness by evaluating a number of different network architectures. The
547 following problems arise in such algorithms. (i) Curse of dimensionality: For each additional
548 hyper-parameter, the size of the search space grows exponentially. (ii) Training time: De-
549 pending on the size of the model, the size of the training dataset, and the available hardware,
550 the evaluation of a network architecture might take a considerable amount of computation
551 time. (iii) Interactions between hyper-parameters: It is not enough to consider each hyper-
552 parameter separately, because the best choice for some hyper-parameter might depend on
553 the chosen configurations of other hyper-parameters. (iv) Non-deterministic model fitness:
554 The fitting of a NN is a non-deterministic process and the weights of a model might not
555 converge to the same optimum in repeated runs. This means that a single evaluation of
556 a network architecture might not reflect the actual fitness of the architecture in general.
557 For an introduction to hyper-parameter optimization in general and the concepts mentioned
558 above in particular, see Hutter et al. (2019). To our knowledge, the following algorithm has
559 not been proposed before.

560 To find a satisfactory model architecture despite the problems listed above, the proposed
561 algorithm works according to the principle of Exploration & Exploitation, which is also
562 explained in Hutter et al. (2019): At the beginning of the search, architectures across the
563 whole search space are evaluated. With an increasing number of evaluations, promising
564 candidates are prioritized.

565 In the following, we consider the search space $H = H_1 \times \dots \times H_n$ being the Cartesian
566 product of a family of domains $H_1; \dots; H_n$ of n hyper-parameters for some integer $n \geq 1$.
567 The set H_i consists of $m_i \geq 1$ available configurations of the i -th hyper-parameter, i.e.,
568 $H_i = \{c_i^1; \dots; c_i^{m_i}\}$ for each $i = 1; \dots; n$.

569 *a. Fitness of an evaluation*

570 In each iteration of the hyper-parameter optimization algorithm, the fitness $f(h)$ of a
571 model architecture specified by $h = (c_1; \dots; c_n) \in H$ is evaluated, where $c_i \in H_i$ for each
572 $i = 1; \dots; n$. The model architectures considered in this paper were trained for 6 epochs on a
573 training data set (April + May 2016) and validated after each epoch on a separate validation
574 data set (June 2016). Note that an epoch refers to one pass-through of the training dataset
575 in the training process. We define the fitness $f(h)$ of a configuration $h \in H$ as the smallest
576 model error achieved in any of the 6 epochs, whereas the model error is determined by the
577 loss function (Chollet 2017) of the NN. Since we consider a classification problem in this
578 paper, the loss function “categorical cross entropy“ is used (Alla and Adari 2019).

579 To find out which number of epochs is sufficient, the model errors for 20 epochs have been
580 determined for a number of model architectures. However, for most model architectures
581 the minimum model error converged within the first 5 epochs. Therefore, we considered 6
582 epochs in order to derive the results obtained in this paper.

583 *b. Selection of new hyper-parameter configurations*

584 A common strategy to pick model configurations for evaluation is the so-called random
585 search method, where a certain probability distribution, e.g. the uniform distribution, is

586 considered on the search space H from which new model architectures are sampled (Hutter
 587 et al. 2019).

588 The idea of the algorithm presented in this section is to start with a random search.
 589 However, after having made a number of evaluations, we can already estimate how the
 590 configuration $c_i \in H_i$ of a single hyper-parameter, for some $i \in \{1, \dots, n\}$, affects the
 591 fitness of the model architecture $h = (c_1, \dots, c_n)$. Based on this information, the probability
 592 distribution on H , from which further model architectures are sampled, can be adapted
 593 to favor model architectures which are more likely to perform well. With an increasing
 594 number of evaluations, the same concept can be applied to an increasing number of j hyper-
 595 parameters, where $j \in \{2, \dots, n\}$, in order to find out which configurations $h' \in H'_j =$
 596 $H_{i_1} \times \dots \times H_{i_j}$ for some subset $J = \{i_1, \dots, i_j\} \subseteq \{1, \dots, n\}$ perform well in combination
 597 with each other. In the following, we sometimes write H'_{h^0} instead of H'_J , in cases where we
 598 want to emphasize that a specific partial architecture h' belongs to the set H'_J . Furthermore,
 599 let $H^* = \bigcup_{J \in \mathcal{P}(\{1, \dots, n\})} H'_J$ denote the set of all (partial) model architectures.

600 1) Definitions

The $(k + 1)$ -th choice of the model configuration $h_{k+1} \in H$, which is to be evaluated
 next, is made based on the set of previously evaluated configurations $E = \{h_1, \dots, h_k\}$, for
 which the fitnesses $f(h_i); i = 1, \dots, k$ have been determined as described in Section 4a.
 Let $h' \in H'_J$ be a partial model architecture for a subset of hyper-parameters with indices
 $J = \{i_1, \dots, i_j\} \subseteq \{1, \dots, n\}$. Then define

$$E_{h^0} = \{h \in E : h' \subseteq h\}$$

as the set of all evaluated hyper-parameter configurations h that share the same configurations with the partial architecture \mathcal{H} . For a given integer $s \in \mathbb{N} = \{1; 2; \dots; g\}$ we define

$$E_s = \{ \mathcal{H} \subseteq \mathcal{H}^* : |E_{\mathcal{H}}| \geq s; |E_{\mathcal{H} \cup \{c_i\}}| < s \text{ for all } c_i \in \mathcal{H}_i; i \in \{1; \dots; n\} \}$$

601 which is the set of partial architectures \mathcal{H} with at least s evaluations and for which all
 602 extensions $\mathcal{H} \cup \{c_i\}$ have less evaluations than s . In other words, E_s contains the largest
 603 partial architectures for which a minimum number s of evaluations exist.

When considering a partial architecture $\mathcal{H} \subseteq \mathcal{H}'$ for evaluation, we are not only interested whether E contains enough data to estimate $f(\mathcal{H})$, but also if all partial architectures in \mathcal{H}'_{h^0} have enough evaluations. Hence we define

$$E_s^p = \{ \mathcal{H} \subseteq E_s : |E_{\mathcal{H}}| \geq \min_{g^0 \in \mathcal{H}'_{h^0}} |E_{g^0}|^p \}$$

604 for a given value $p \in [1; 1)$, i.e., E_s^p is a subset of E_s containing all partial architectures
 605 \mathcal{H} with a number of evaluations $|E_{\mathcal{H}}|$ below an upper bound, which depends on the partial
 606 architecture $g^0 \in \mathcal{H}'_{h^0}$ with the smallest number of evaluations $|E_{g^0}|$.

Furthermore, we define the median fitness $M_E(\mathcal{H})$ for a partial model architecture \mathcal{H} as

$$M_E(\mathcal{H}) = \text{median}(f(h) : h \in E_{\mathcal{H}})$$

and, finally, the set $\Delta(\mathcal{H})$ of partial architectures g' , which share $|g' \cap \mathcal{H}|$ configurations with \mathcal{H} as

$$\Delta(\mathcal{H}) = \{ g' \subseteq \mathcal{H} : |g' \cap \mathcal{H}| = |g'| \text{ for } g' \in \mathcal{H} \}$$

607 2) The algorithm

608 In this section we explain how we determine an initial partial architecture \mathcal{H}_0 and how we
 609 pick subsequent partial architectures $\mathcal{H}_1; \dots; \mathcal{H}_k$ until their union is a full architecture that
 610 can be evaluated next and added to the set E .

For given values of s and ρ , we sample from the set of partial architectures E_s^ρ . Initially, the set E_s^ρ is empty, because E is empty since no architectures have been evaluated, yet. Note that E_s^ρ can also be empty due to the upper bound determined by the parameter ρ . In these cases, a partial architecture $h'_0 = f c_i g$ with a random configuration $c_i \geq H_i$ for a random hyper-parameter H_i is chosen. Otherwise, we sample h'_0 from the set E_s^ρ with probability $P_{E_s^\rho}$, where P_F is defined as

$$P_F(\mathcal{H}) = \mathbb{P} \frac{2^{-M_E(\mathcal{H})=d}}{g^0 \in F 2^{-M_E(g^0)=d}}$$

611 for any partial architecture $\mathcal{H} \geq F$, a given set F of partial architectures and some $d > 0$,
612 which is another parameter of the algorithm, along with s and ρ . Note that P_F is defined such
613 that a partial architecture \mathcal{H} is twice as likely to be chosen than \mathcal{G} if $M_E(\mathcal{H}) = M_E(\mathcal{G}) + d$.
614 Furthermore, once enough evaluations have been made such that E_s^ρ is non-empty, the
615 algorithm will always pick a partial architecture from E_s^ρ . Without the upper bound, the
616 first few partial architectures in E_s would be sampled ad infinitum, while other partial
617 architectures, which have not been sampled often enough yet, would not be sampled at all.
618 Therefore it is necessary to include the upper bound in E_s^ρ for the number of evaluations.

619 We now have the first part h'_0 of the architecture h , which we want to evaluate. Next,
620 we successively determine more parts $h'_1; \dots; h'_k$ until their union is a complete architecture
621 $h = h'_0 [\dots [h'_k \geq H$. For this, let $\bar{h}_j = h'_0 [\dots [h'_j$ be the union of all partial architectures
622 up to h'_j .

623 For $j \geq 1; \dots; k$ we iteratively sample h'_j from $E_s^\rho \setminus \Delta(\bar{h}_{j-1})$ with probability
624 $P_{E_s^\rho \setminus \Delta(\bar{h}_{j-1})}(h'_j)$ with the smallest $\geq \mathbb{N}$ for which $E_s^\rho \setminus \Delta(\bar{h}_{j-1})$ is not empty. In other
625 words, we choose h'_j such that it has at least one new configuration and also the largest
626 possible overlap with \bar{h}_{j-1} , and that the new configurations are likely to perform well in

627 combination with the previously chosen configurations. If $E_s^p \setminus \Delta(\bar{h}_{j-1})$ is empty for all
628 $2 \leq j \leq N$, the partial configuration $h_j = f c_j g$ consists of a single randomly chosen configuration
629 $c_j \in H_j$, where H_j is a random hyper-parameter, for which it holds that $H_j \setminus \bar{h}_{j-1} = \emptyset$.
630 Once enough architectures have been evaluated and we want to pick the best architecture
631 based on E , we follow the same steps as described above, but instead of sampling partial
632 architectures from E_s^p with probability $P_{E_s^p}$, we pick the partial architectures $h' \in E_s$ with
633 the lowest $M_E(h')$ instead.

634 References

- 635 Alla, S. and Adari, S. K. (2019). *Beginning Anomaly Detection Using Python-Based Deep*
636 *Learning*. Springer. ISBN: 978-1484251775.
- 637 Armstrong, J. S. and Grohman, M. C. (1972). A comparative study of methods for long-
638 range market forecasting. *Management Science*, 19:211–221.
- 639 Ben Bouallègue, Z. (2013). Calibrated short-range ensemble precipitation forecasts using
640 extended logistic regression with interaction terms. *Weather and Forecasting*, 28:515–524.
- 641 Bergstra, J. S., Bardenet, R., Bengio, Y., and Kégl, B. (2011). Algorithms for hyper-
642 parameter optimization. In Shawe-Taylor, J., Zemel, R. S., Bartlett, P. L., Pereira, F.,
643 and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems*, pages
644 2546–2554. Curran Associates, Inc.
- 645 Bouttier, F. and Marchal, H. (2020). Probabilistic thunderstorm forecasting by blending
646 multiple ensembles. *Tellus A*, 72:1–19.

647 Bowler, N. E., Pierce, C. E., and Seed, A. W. (2006). STEPS: A probabilistic precipita-
648 tion forecasting scheme which merges an extrapolation nowcast with downscaled NWP.
649 *Quarterly Journal of the Royal Meteorological Society*, 132:2127–2155.

650 Bridle, J. S. (1990). Probabilistic interpretation of feedforward classification network out-
651 puts, with relationships to statistical pattern recognition. In Soulié, F. F. and Héroult,
652 J., editors, *Neurocomputing*, pages 227–236. Springer.

653 Brunet, G., Jones, S., and Ruti, P. M. (2015). *Seamless Prediction of the Earth System:
654 from Minutes to Months*. World Meteorological Organization. ISBN: 978-9263111562.

655 Chollet, F. (2017). *Deep Learning with Python*. Manning Publications. ISBN: 978-
656 1617294433.

657 Dozat, T. (2016). Incorporating nesterov momentum into adam. In *International Conference
658 on Learning Representations, ICLR 2016, Workshop Track*. [https://openreview.net/
659 pdf/OM0jvwB8jlp57ZJjtNEZ.pdf](https://openreview.net/pdf/OM0jvwB8jlp57ZJjtNEZ.pdf).

660 Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive subgradient methods for online
661 learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–
662 2159.

663 Foresti, L. and Seed, A. (2014). The effect of flow and orography on the spatial distribution
664 of the very short-term predictability of rainfall from composite radar images. *Hydrology
665 and Earth System Sciences*, 18:4671.

666 Gebhardt, C., Theis, S., Paulat, M., and Bouallègue, Z. B. (2011). Uncertainties in cosmo-de
667 precipitation forecasts introduced by model perturbations and variation of lateral bound-
668 aries. *Atmospheric Research*, 100:168–177.

- 669 Germann, U. and Zawadzki, I. (2002). Scale-dependence of the predictability of precipitation
670 from continental radar images. Part I: Description of the methodology. *Monthly Weather*
671 *Review*, 130:2859–2873.
- 672 Glorot, X., Bordes, A., and Bengio, Y. (2011). Deep sparse rectifier neural networks. In Gor-
673 don G., Dunson D. and Dudík M., editors, *Proceedings of the Fourteenth International*
674 *Conference on Artificial Intelligence and Statistics, JMLR Workshop and Conference Pro-*
675 *ceedings*. pages 15:315–323.
- 676 Gneiting, T., Raftery, A. E., Westveld III, A. H., and Goldman, T. (2005). Calibrated
677 probabilistic forecasting using ensemble model output statistics and minimum CRPS es-
678 timation. *Monthly Weather Review*, 133:1098–1118.
- 679 Golding, B. (1998). Nimrod: A system for generating automated very short range forecasts.
680 *Meteorological Applications*, 5:1–16.
- 681 Haiden, T., Kann, A., Wittmann, C., Pistotnik, G., Bica, B., and Gruber, C. (2011). The
682 Integrated Nowcasting through Comprehensive Analysis (INCA) system and its validation
683 over the Eastern Alpine region. *Weather and Forecasting*, 26:166–183.
- 684 Hamill, T. M., Hagedorn, R., and Whitaker, J. S. (2008). Probabilistic forecast calibration
685 using ECMWF and GFS ensemble reforecasts. Part II: Precipitation. *Monthly Weather*
686 *Review*, 136:2620–2632.
- 687 Hess, R. (2020). Statistical postprocessing of ensemble forecasts for severe weather at
688 Deutscher Wetterdienst. *Nonlinear Processes in Geophysics*. DOI: 10.5194/npg-2019-64
689 (to appear).

- 690 Hutter, F., Kotthoff, L., and Vanschoren, J., editors (2019). *Automated Machine Learning:
691 Methods, Systems, Challenges*. Springer.
- 692 Johnson, A. and Wang, X. (2012). Verification and calibration of neighborhood and object-
693 based probabilistic precipitation forecasts from a multimodel convection-allowing ensem-
694 ble. *Monthly Weather Review*, 140:3054–3077.
- 695 Keras (2020). Keras API reference: backend utilities. [https://keras.io/api/uti ls/
696 backend_uti ls/](https://keras.io/api/uti ls/backend_uti ls/). accessed on 2020.08.20.
- 697 Kingma, D. P. and Ba, J. (2015). Adam: A method for stochastic optimization. In Bengio,
698 Y. and LeCun, Y., editors, *3rd International Conference on Learning Representations,
699 ICLR 2015, Conference Track Proceedings*. <http://arxiv.org/abs/1412.6980>.
- 700 Kober, K., Craig, G., Keil, C., and Dörnbrack, A. (2012). Blending a probabilistic nowcast-
701 ing method with a high-resolution numerical weather prediction ensemble for convective
702 precipitation forecasts. *Quarterly Journal of the Royal Meteorological Society*, 138:755–
703 768.
- 704 Murphy, A. H. and Winkler, R. L. (1977). Reliability of subjective probability forecasts of
705 precipitation and temperature. *Journal of the Royal Statistical Society: Series C*, 26:41–
706 47.
- 707 Murphy, A. H. and Winkler, R. L. (1987). A general framework for forecast verification.
708 *Monthly Weather Review*, 115:1330–1338.
- 709 Nerini, D., Foresti, L., Leuenberger, D., Robert, S., and Germann, U. (2019). A reduced-
710 space ensemble kalman filter approach for flow-dependent integration of radar extrapola-
711 tion nowcasts and nwp precipitation ensembles. *Monthly Weather Review*, 147:987–1006.

712 Nicolis, C., Perdigao, R. A., and Vannitsem, S. (2009). Dynamics of prediction errors under
713 the combined effect of initial condition and model errors. *Journal of the Atmospheric*
714 *Sciences*, 66:766–778.

715 Palmer, T. N. (2002). The economic value of ensemble forecasts as a tool for risk assessment:
716 From days to decades. *Quarterly Journal of the Royal Meteorological Society*, 128:747–774.

717 Park, J. and Sandberg, I. W. (1991). Universal approximation using radial-basis-function
718 networks. *Neural Computation*, 3:246–257.

719 Schaumann, P., de Langlard, M., Hess, R., James, P., and Schmidt, V. (2020). A calibrated
720 combination of probabilistic precipitation forecasts to achieve a seamless transition from
721 nowcasting to very short-range forecasting. *Weather and Forecasting*, 35:773–791.

722 Scheuerer, M. (2014). Probabilistic quantitative precipitation forecasting using ensemble
723 model output statistics. *Quarterly Journal of the Royal Meteorological Society*, 140:1086–
724 1096.

725 Schwartz, C. S. and Sobash, R. A. (2017). Generating probabilistic forecasts from convection-
726 allowing ensembles using neighborhood approaches: A review and recommendations.
727 *Monthly Weather Review*, 145:3397–3418.

728 Seed, A. W., Pierce, C. E., and Norman, K. (2013). Formulation and evaluation of a scale
729 decomposition-based stochastic precipitation nowcast scheme. *Water Resources Research*,
730 49:6624–6641.

731 Stephan, K., Klink, S., and Schraff, C. (2008). Assimilation of radar-derived rain rates
732 into the convective-scale model COSMO-DE at DWD. *Quarterly Journal of the Royal*
733 *Meteorological Society*, 134:1315–1326.

734 Sutskever, I., Martens, J., Dahl, G., and Hinton, G. (2013). On the importance of initializa-
735 tion and momentum in deep learning. Proceedings of the 30th International Conference on
736 Machine Learning, Proceedings of Machine Learning Research (PMLR), pages 28:1139–
737 1147.

738 Theis, S., Gebhardt, C., and Bouallegue, Z. B. (2012). *Beschreibung des COSMO-*
739 *DE-EPS und seiner Ausgabe in die Datenbanken des DWD*. Deutscher Wetterdi-
740 enst. [https://www.dwd.de/SharedDocs/downloads/DE/modelldokumentationen/nwv/](https://www.dwd.de/SharedDocs/downloads/DE/modelldokumentationen/nwv/cosmo_de_eps/cosmo_de_eps_dbbeschr_201208.pdf)
741 [cosmo_de_eps/cosmo_de_eps_dbbeschr_201208.pdf](https://www.dwd.de/SharedDocs/downloads/DE/modelldokumentationen/nwv/cosmo_de_eps/cosmo_de_eps_dbbeschr_201208.pdf).

742 Theis, S., Hense, A., and Damrath, U. (2005). Probabilistic precipitation forecasts from a
743 deterministic model: A pragmatic approach. *Meteorological Applications*, 12:257–268.

744 Vannitsem, S., Bremnes, J. B., Demaeyer, J., Evans, G. R., Flowerdew, J., Hemri, S., Lerch,
745 S., Roberts, N., Theis, S., Atencia, A., et al. (2020). Statistical postprocessing for weather
746 forecasts—review, challenges and avenues in a big data world. Preprint: *arXiv:2004.06582*.

747 Venugopal, V., Foufoula-Georgiou, E., and Sapozhnikov, V. (1999). Evidence of dynamic
748 scaling in space-time rainfall. *Journal of Geophysical Research*, 104:31599–31610.

749 Wang, Y., Coning, E., Harou, A., Jacobs, W., Joe, P., Nikitina, L., Roberts, R., Wang,
750 J., and Wilson, J. (2017). Guidelines for nowcasting techniques. *WMO Publications*.
751 published online: https://library.wmo.int/doc_num.php?expl_num_id=3795.

752 Weigl, E. and Winterrath, T. (2010). Radargestützte Niederschlagsanalyse und -vorhersage
753 (RADOLAN, RADVOR-OP). *Promet*, 35:78–86.

754 Wilks, D. S. (2006). *Statistical Methods in the Atmospheric Sciences*, volume 91. Academic
755 Press.

- 756 Wilks, D. S. (2009). Extending logistic regression to provide full-probability-distribution
757 mos forecasts. *Meteorological Applications*, 16:361–368.
- 758 Winterrath, T. and Rosenow, W. (2007). A new module for the tracking of radar-derived
759 precipitation with model-derived winds. *Advances in Geosciences*, 10:77–83.
- 760 Winterrath, T., Rosenow, W., and Weigl, E. (2012). On the DWD quantitative precipitation
761 analysis and nowcasting system for real-time application in German flood risk manage-
762 ment. In Moore, R. J., Cole, S. J., and Illingworth, A. J., editors, *Weather Radar and*
763 *Hydrology*, IAHS Proceedings and Reports, pages 351:323–329.
- 764 Zawadzki, I., Morneau, J., and Laprise, R. (1994). Predictability of precipitation patterns:
765 An operational approach. *Journal of Applied Meteorology*, 33:1562–1571.
- 766 Zeiler, M. D. (2012). Adadelta: an adaptive learning rate method. Preprint:
767 *arXiv:1212.5701*.
- 768 Zhang, W., Itoh, K., Tanida, J., and Ichioka, Y. (1990). Parallel distributed processing
769 model with local space-invariant interconnections and its optical architecture. *Applied*
770 *Optics*, 29:4790–4797.

771 **LIST OF TABLES**

772 **Table 1.** Selected configurations of hyper-parameters for different lead times . . . 40

Lead time	Nr. of conv. layers	Kernel size	Conv. activation	Conv. reg.	Conv. Output length
+1h	4	3	elu	5e-07	14
+2h	4	3	elu	5e-06	6
+3h	5	3	relu	0	8
+4h	1	9	elu	0	4
+5h	1	9	elu	5e-06	12
+6h	2	6	sigmoid	0	4

Lead time	Nr. of Neurons in dense layer	Dense activation	Dense reg.	Nr. of triangular func.	Optimizer
+1h	12	sigmoid	0	9	Nadam
+2h	12	tanh	0	5	Nadam
+3h	12	exponential	1e-06	3	Adamax
+4h	6	tanh	0	5	Adamax
+5h	10	sigmoid	0	3	Adamax
+6h	10	relu	1e-05	5	Adamax

Table 1. Selected configurations of hyper-parameters for different lead times

773 **LIST OF FIGURES**

774 **Fig. 1.** Validation scores for 5 of the 9 considered thresholds for the NN-model (green),
 775 Ensemble-MOS (yellow) and RadVOR (blue). Left: bias, middle: Brier skill
 776 score, right: reliability. 42

777 **Fig. 2.** Reliability diagrams for the NN-model (green), Ensemble-MOS (yellow) and Rad-
 778 VOR (blue) for 5 of the 9 considered thresholds (columns) and 6 lead times (rows).
 779 Bins with less than 100 exceedance probabilities are omitted. 43

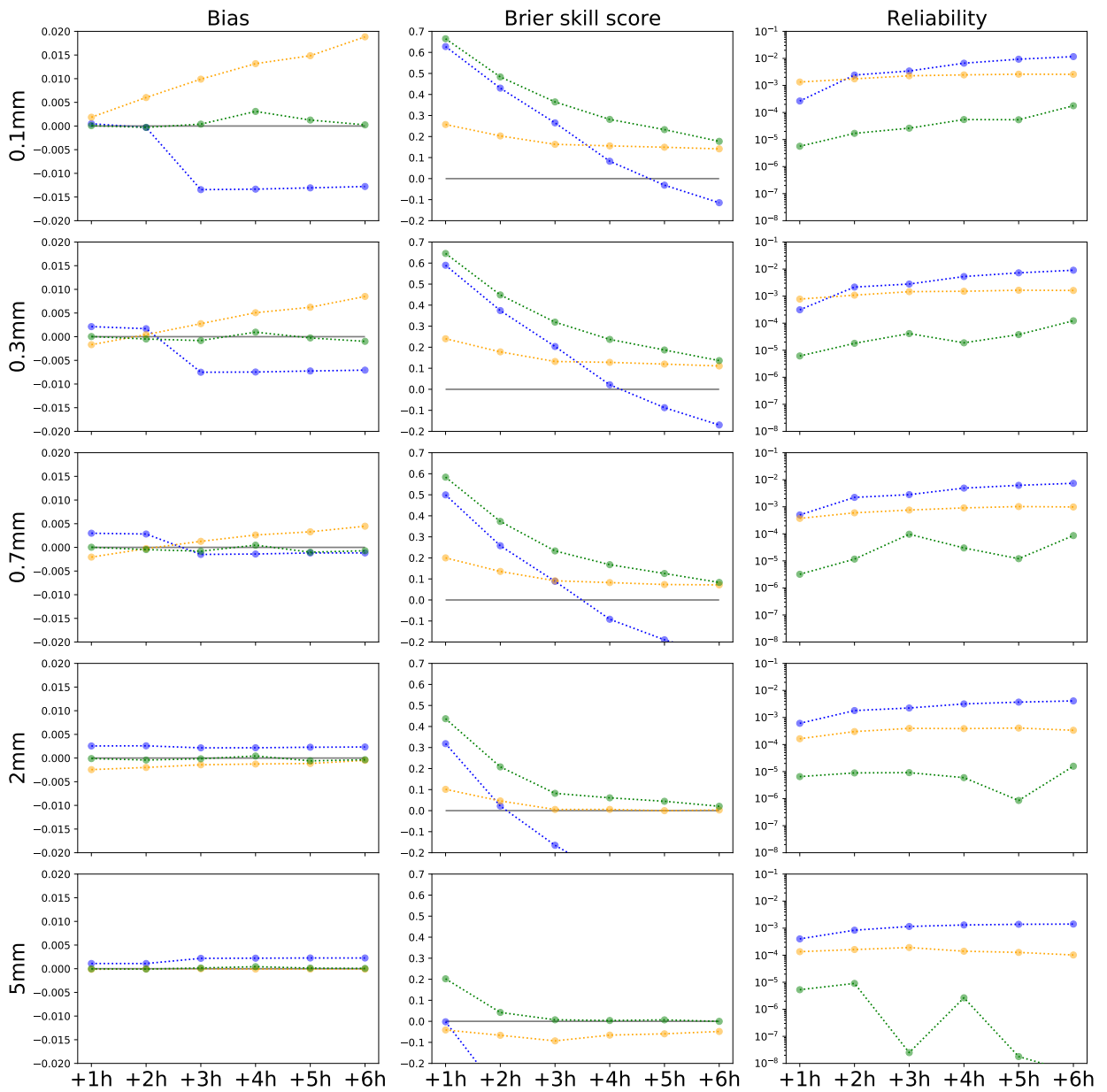
780 **Fig. 3.** A schematic representation of the network architecture (green) and the input and
 781 output data (blue). The arrows indicate the flow of information. 44

782 **Fig. 4.** Distribution of model fitness (categorical cross entropy) for the considered con-
 783 figurations of the first 9 hyper-parameters and the 5 optimizers, for lead time
 784 +1 h (left) and +6 h (right). The median fitness is depicted by a blue line. Due
 785 to very long tails, the distributions are only shown between the respective 5th
 786 and 95th percentiles. 45

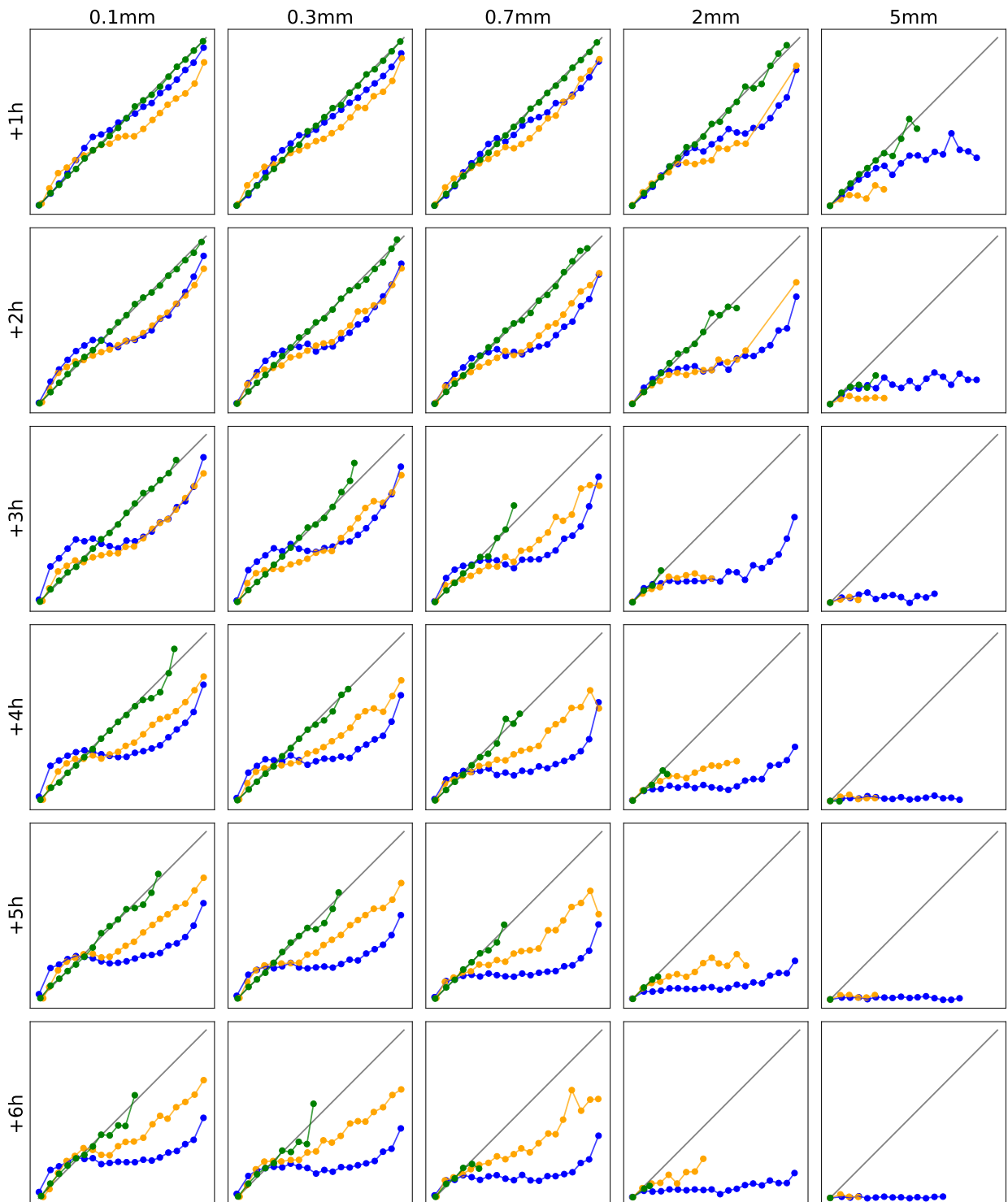
787 **Fig. 5.** Median model fitness (categorical cross entropy) for different hyper-parameter
 788 settings regarding the convolutional layers, and for several lead times: +1 h,
 789 +3 h, +5 h, +6 h (top to bottom). The color of each tile depicts the median
 790 model fitness in dependence on number and kernel size of the convolutional layers
 791 (*x*-axis) and activation function (*y*-axis). Whereby the *x*-axis labels correspond
 792 to the kernel size of each layer, e.g., 5 5 5 stands for three layers with a kernel
 793 size of five each. The number in each tile indicates the number of evaluations
 794 made by the hyper-parameter optimization algorithm. 46

795 **Fig. 6.** Median model fitness (categorical cross entropy) for different hyper-parameter
 796 settings regarding triangular functions (*x*-axis) and interaction terms (*y*-axis),
 797 and for the lead times +1 h (top-left), +2 h (top-right), +3 h (middle-left), +4 h
 798 (middle-right), +5 h (bottom-left) and +6 h (bottom-right). The color depicts
 799 the median model fitness. The number in each tile indicates the number of
 800 evaluations made by the hyper-parameter optimization algorithm. 47

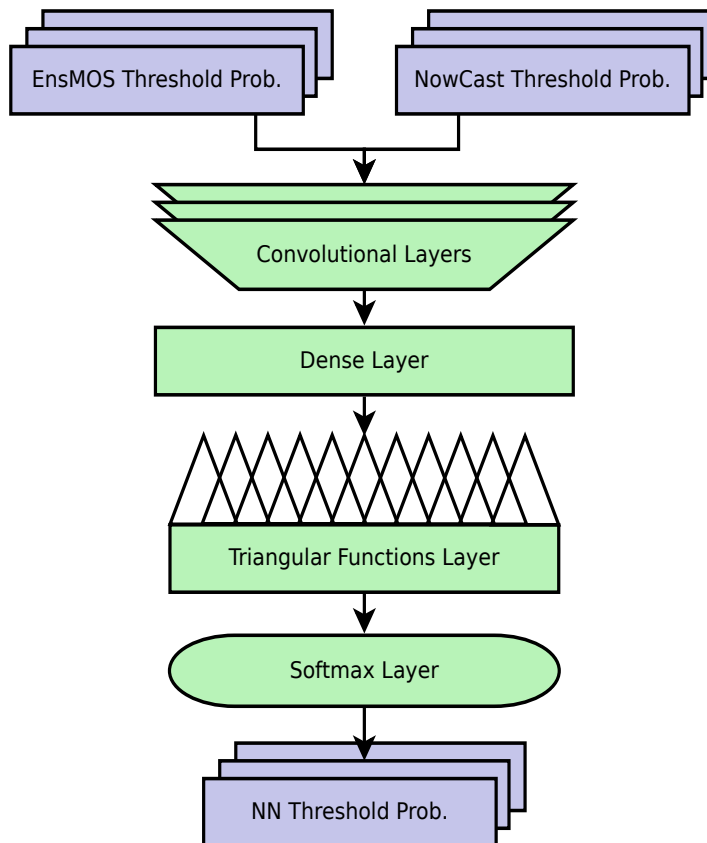
801 **Fig. 7.** Combination example for +1 h (first row), +2 h (second row) and +3 h (third
 802 row), for the hour 6-7 am of 2016-07-21. Exceedance probabilities for each thresh-
 803 old are illustrated by the area of circles and corresponding color. The area of grey
 804 circles corresponds to probability one. The grey solid lines indicate the borders
 805 of federal states of Germany with Hesse and Saxony-Anhalt in the centre. 48



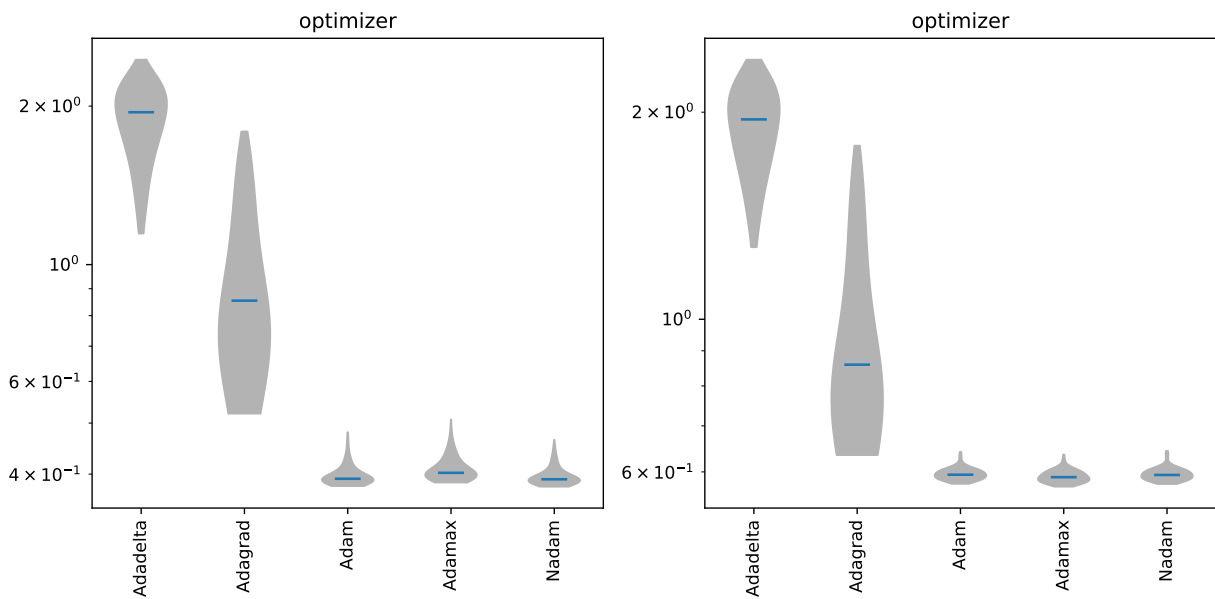
806 Fig. 1. Validation scores for 5 of the 9 considered thresholds for the NN-model (green), Ensemble-
 807 MOS (yellow) and RadVOR (blue). Left: bias, middle: Brier skill score, right: reliability.



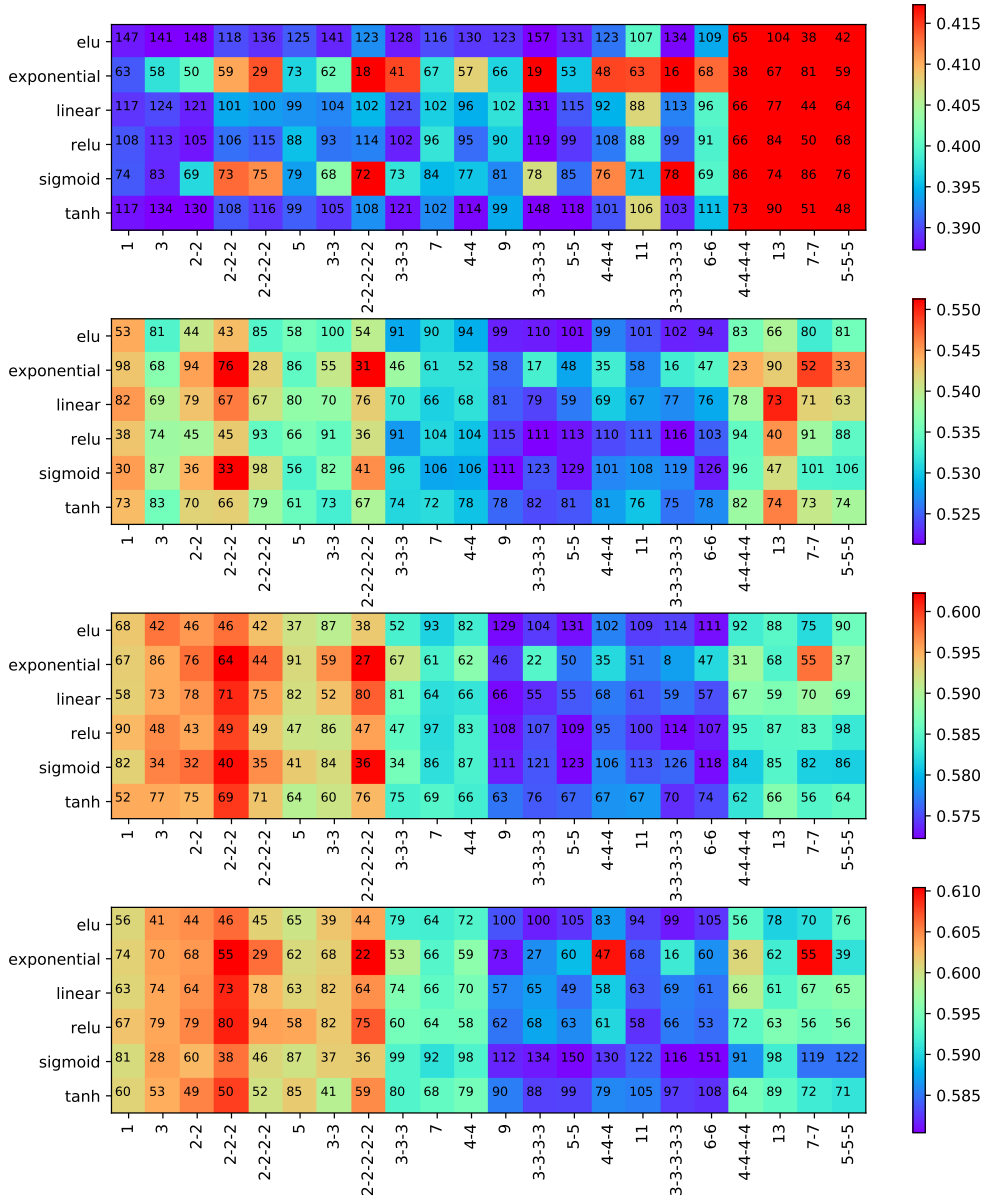
808 Fig. 2. Reliability diagrams for the NN-model (green), Ensemble-MOS (yellow) and RadVOR
 809 (blue) for 5 of the 9 considered thresholds (columns) and 6 lead times (rows). Bins with less than
 810 100 exceedance probabilities are omitted.



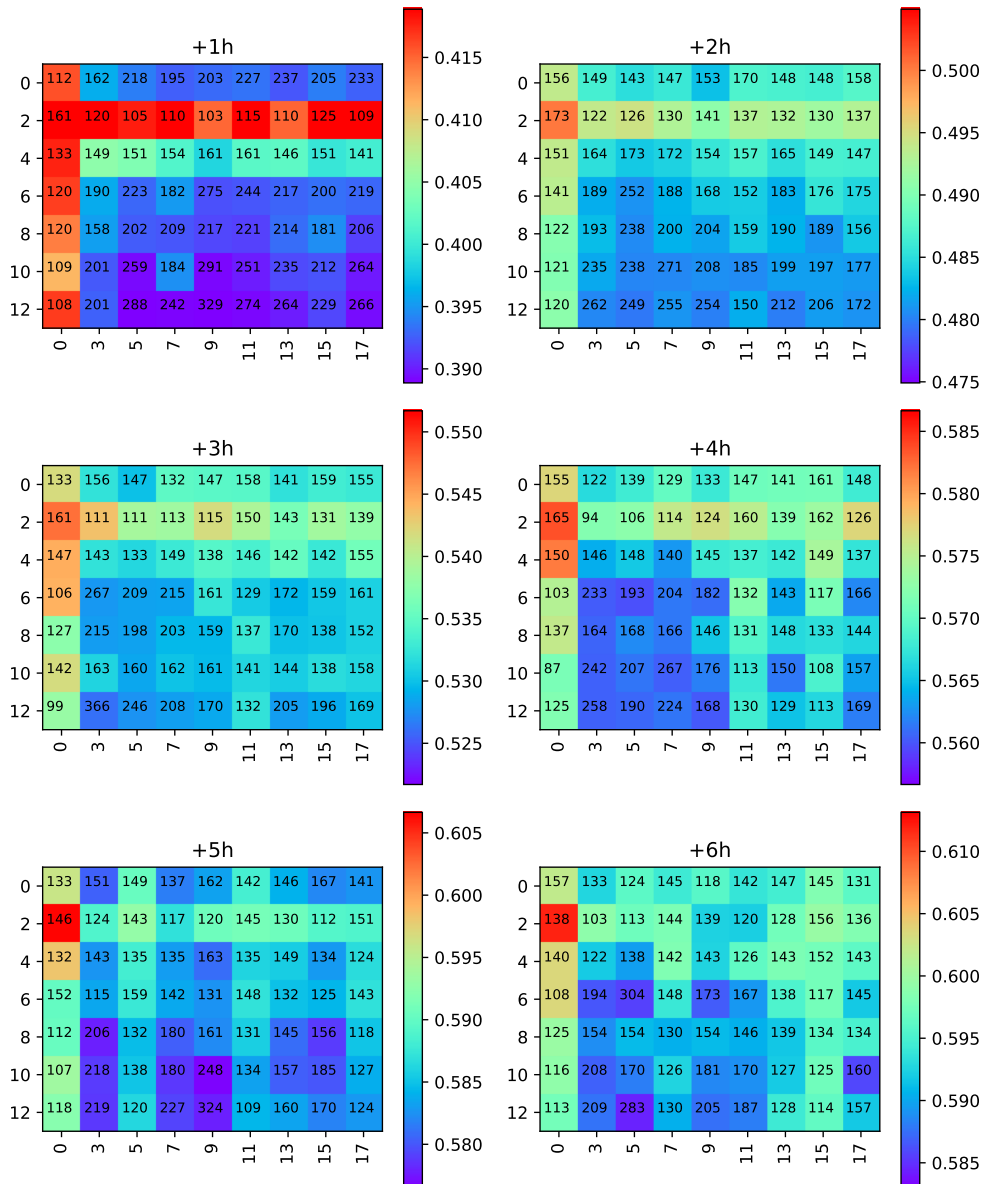
811 Fig. 3. A schematic representation of the network architecture (green) and the input and output
 812 data (blue). The arrows indicate the flow of information.



813 Fig. 4. Distribution of model fitness (categorical cross entropy) for the considered configurations
 814 of the first 9 hyper-parameters and the 5 optimizers, for lead time +1 h (left) and +6 h (right).
 815 The median fitness is depicted by a blue line. Due to very long tails, the distributions are only
 816 shown between the respective 5th and 95th percentiles.



817 Fig. 5. Median model fitness (categorical cross entropy) for different hyper-parameter settings
 818 regarding the convolutional layers, and for several lead times: +1h, +3h, +5h, +6h (top to
 819 bottom). The color of each tile depicts the median model fitness in dependence on number and
 820 kernel size of the convolutional layers (x -axis) and activation function (y -axis). Whereby the x -
 821 axis labels correspond to the kernel size of each layer, e.g., 5 5 5 stands for three layers with a
 822 kernel size of five each. The number in each tile indicates the number of evaluations made by the
 823 hyper-parameter optimization algorithm.



824 Fig. 6. Median model fitness (categorical cross entropy) for different hyper-parameter settings
 825 regarding triangular functions (x -axis) and interaction terms (y -axis), and for the lead times +1 h
 826 (top-left), +2 h (top-right), +3 h (middle-left), +4 h (middle-right), +5 h (bottom-left) and +6 h
 827 (bottom-right). The color depicts the median model fitness. The number in each tile indicates the
 828 number of evaluations made by the hyper-parameter optimization algorithm.

