

1 **Advances in Land Surface Model-based Forecasting: A Comparison of LSTM,**
2 **Gradient Boosting, and Feedforward Neural Networks as Prognostic State Emulators in**
3 **a Case Study with ECLand**

4
5 Marieke Wesselkamp¹, Matthew Chantry², Ewan Pinnington², Margarita Choulga², Souhail
6 Boussetta², Maria Kalweit³, Joschka Boedecker^{3,4}, Carsten F. Dormann¹, Florian
7 Pappenberger², and Gianpaolo Balsamo^{2,5}

8
9
10 1 Department of Biometry, University of Freiburg, Germany

11 2 European Centre for Medium-Range Weather Forecasts, Reading, United Kingdom

12 3 Department of Computer Science, University of Freiburg, Germany

13 4 BrainLinks-BrainTools, University of Freiburg, Germany

14 5 World Meteorological Organization, Geneva, Switzerland

15

16

17 Correspondence to: Marieke Wesselkamp (marieke.wesselkamp@biom.uni-freiburg.de)

18

Abstract

19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

Most useful weather prediction for the public is near the surface. The processes that are most relevant for near-surface weather prediction are also those that are most interactive and exhibit positive feedback or have key roles in energy partitioning. Land surface models (LSMs) consider these processes together with surface heterogeneity and, when coupled with an atmospheric model, provide boundary and initial conditions. They forecast water, carbon and energy fluxes, which are an integral component of coupled atmospheric models. This numerical parametrization of atmospheric boundaries is computationally expensive and statistical surrogate models are increasingly used to accelerate experimental research. We evaluated the efficiency of three surrogate models in simulating land surface processes for speeding up experimental research. Specifically, we compared the performance of a Long-Short Term Memory (LSTM) encoder-decoder network, extreme gradient boosting, and a feed-forward neural network within a physics-informed multi-objective framework. This framework emulates key prognostic states of the ECMWF's Integrated Forecasting System (IFS) land surface scheme, ECLand, across continental and global scales. Our findings indicate that while all models on average demonstrate high accuracy over the forecast period, the LSTM network excels in continental long-range predictions when carefully tuned, XGB scores consistently high across tasks and the MLP provides an excellent implementation-time-accuracy trade-off. While their reliability is context dependent, the runtime reductions achieved by the emulators in comparison to the full numerical models are significant, offering a faster alternative for conducting experiments on land surfaces.

42 **1 Introduction**

43

44 While forecasting of climate and weather system processes has long been a task for numerical
45 models, recent developments in deep learning have introduced competitive machine-learning
46 (ML) systems for numerical weather prediction (NWP) (Bi et al., 2022; Lam et al., 2023;
47 Lang et al., 2024). Land surface models (LSMs), even though being an integral part of
48 numerical weather prediction, have not yet caught the attention of the ML-community. LSMs
49 forecast water, carbon and energy fluxes and, in coupling with an atmospheric model, provide
50 the lower boundary and initial conditions (Boussetta et al., 2021; De Rosnay et al.,
51 2014). The parametrization of land surface states does not only affect predictability of earth
52 and climate systems on sub-seasonal scales (Muñoz-Sabater et al., 2021), but also the short-
53 and medium-range skill of NWP forecasts (De Rosnay et al., 2014). Beyond their online
54 integration with NWPs, offline versions of LSMs provide research tools for experiments on
55 the land surface (Boussetta et al., 2021), the diversity of which, however, are limited by
56 substantial computational resources requirements and often moderate runtime efficiencies
57 (Reichstein et al., 2019).

58 Emulators constitute statistical surrogates for numerical simulation models that, by
59 approximating the latter, aim for increasing computational efficiency (Machac et al., 2016).
60 While the construction of emulators can itself require substantial computational resources,
61 their subsequent evaluation usually runs orders of magnitude faster than the original
62 numerical model (Fer et al., 2018). For this reason, emulators have found application for
63 example in modular parametrization of online weather forecasting systems (Chantry et al.,
64 2021), in replacing the MCMC-sampling procedure in Bayesian calibration of ecosystem
65 models (Fer et al., 2018), or in generating forecast ensembles of atmospheric states for
66 uncertainty quantification (Li et al., 2023). Beyond their computational efficiency, surrogate
67 models with high parametric flexibility have the potential to correct process mis-specification
68 in a physical model when fine-tuned to observations (Wesselkamp et al., 2022).

69 Modelling approaches used for emulation range from low parametrized, auto-regressive
70 linear models to highly non-linear and flexible neural networks (Baker et al., 2022; Chantry
71 et al., 2021; Meyer et al., 2022; Nath et al., 2022). In the global land surface system M-
72 MESMER, a set of simple AR1 regression models is used to initialize the numerical LSM,
73 resulting in a modularized emulator (Nath et al., 2022). Numerical forecasts of gross primary
74 productivity and hydrological targets were successfully approximated by Gaussian processes
75 (Baker et al., 2022; Machac et al., 2016), the advantage of which is their direct quantification

76 of prediction uncertainty. When it comes to highly diverse or structured data, neural networks
77 have shown to deliver accurate approximations, for example for gravity wave drags and
78 urban surface temperature (Chantry et al., 2021; Meyer et al., 2022). In most fields of
79 machine learning, specific types of neural networks are now the best approach to representing
80 fit and prediction. One exception is so-called tabular data, i.e. data without spatial or temporal
81 interdependencies (as opposed to vision and sound), where extreme gradient boosting is still
82 the go-to approach (Grinsztajn et al., 2022; Shwartz-Ziv and Armon, 2021).

83 ECLand is the land surface scheme that provides boundary and initial conditions for the
84 Integrated Forecasting System (IFS) of the European Centre for Medium-range Weather
85 Forecasts (ECMWF) (Boussetta et al., 2021). Driven by meteorological forcing and spatial
86 climate fields, it has a strong influence on the NWP (De Rosnay et al., 2014) and also
87 constitutes a standalone framework for offline forecasting of land surface processes (Muñoz-
88 Sabater et al., 2021). The modular construction of ECLand offers potential for element-wise
89 improvement of process representation and thus a stepwise development towards increased
90 computational efficiency. Within the IFS, ECLand also forms the basis of the land surface
91 data assimilation system, updating the land surface state with synoptic data and satellite
92 observations of soil moisture and snow. Emulators of physical systems have been shown to
93 be beneficial in data assimilation routines, allowing for a quick estimation and low
94 maintenance of the tangent linear model (Hatfield et al., 2021). Together with the potential to
95 run large ensembles of land surface states at a much-reduced cost, this would be a potential
96 application of the surrogate models introduced here.

97 Long-short term memory networks (LSTMs) have gained popularity in hydrological
98 forecasting as rainfall-runoff models, for predicting stream flow temperature and also soil
99 moisture (Bassi et al., 2024; Kratzert et al., 2019b; Lees et al., 2022; Zwart et al., 2023).

100 Research on the interpretability of LSTMs has found correlations between the model cell
101 states and spatially or thematically similar hydrological units (Lees et al., 2022), suggesting
102 the specific usefulness of LSTM for representing variables with dynamic storages and
103 reservoirs (Kratzert et al., 2019a). As emulators, LSTMs have been shown useful for sea
104 surface level projection in a variational manner with Monte Carlo dropout (Van Katwyk et al.,
105 2023).

106 While most of these studies trained their models on observations or reanalysis data, our
107 emulator learns the representation from ECLand simulations directly. To our knowledge, a
108 comparison of models without memory mechanisms to an LSTM-based neural network for
109 global land surface emulation has not been conducted before.

110 We emulate seven prognostic state variables of ECLand, which represent core land surface
111 processes: soil water volume and soil temperature, each at three depth layers, and snow cover
112 fraction at the surface layer. The represented variables would allow their coupling to the IFS,
113 yet the emulators do not replace ECLand in its full capabilities. Yet, these three state variables
114 represent the core of the current configuration of ECLand. We specifically focus on the utility
115 of memory mechanisms, highlighting the development of a single LSTM-based encoder-
116 decoder model compared to an extreme gradient boosting approach (XGB) and a multilayer
117 perceptron (MLP), which all perform the same tasks. The LSTM architecture builds on an
118 encoder-decoder network design introduced for flood forecasting (Nearing et al., 2024). To
119 compare forecast skill systematically, the three emulators were compared in long-range
120 forecasting against climatology (Pappenberger et al., 2015). In this work, the emulators are
121 evaluated on ECLand simulations only, i.e. on purely synthetic data, while we anticipate their
122 validation and fine-tuning on observations for future work.

123

124 **2 Methods**

125

126 **2.1 The Land Surface Model: ECLand**

127

128 ECLand is a tiled ECMWF Scheme for surface exchanges over land that represents surface
129 heterogeneity and incorporates land surface hydrology (Balsamo et al., 2011; ECMWF,
130 2017). ECLand computes surface turbulent fluxes of heat, moisture and momentum and skin
131 temperature over different tiles (vegetation, bare soil, snow, interception and water) and then
132 calculates an area-weighted average for the grid-box to couple with the atmosphere
133 (Boussetta et al., 2021). For the overall accuracy of the model, accurate land surface
134 parameterizations are essential (Kimpson et al., 2023) as they e.g. determine the sensible and
135 latent heat fluxes, and provide the lower boundary conditions for enthalpy and moisture
136 equations in the atmosphere (Viterbo, 2002). We emulate three prognostic state variables of
137 ECLand that represent core land surface processes: soil water volume (m^3m^{-3}) and soil
138 temperature (K) at each three depth layers (each at 0 – 7 cm, 7 – 21 cm and 21 – 72 cm) and
139 snow cover fraction (%), aggregated at the surface layer.

140

141 **2.2 Data sources**

142

143 As training data base, global simulation and reanalysis time series from 2010 to 2022 were
144 compiled to *zarr* format at an aggregated 6-hourly temporal resolution. Simulations and
145 climate fields were generated from ECMWFs development cycle CY49R2, ECLand forced
146 by ERA-5 meteorological reanalysis data (Hersbach et al., 2020).

147 There are three main sources of data used for creation of the data base: The first is a selection
148 of surface physiographic fields from ERA5 (Hersbach et al., 2020) and their updated versions
149 (Boussetta et al., 2021; Choulga et al., 2019; Muñoz-Sabater et al., 2021), used as static
150 model input features (X). The second is a selection of atmospheric and surface model fields
151 from ERA5, used as static and dynamic model input features (Y). The third are ECLand
152 simulations, constituting the model’s dynamic prognostic state variables (z) and hence
153 emulator input and target features. A total of 41 static, seasonal and dynamical features were
154 used to create the emulators, see table 1 for an overview of input variables and details on the
155 surface physiographic and atmospheric fields below.

156

157 **2.2.1 Surface physiographic fields**

158

159 Surface physiographic fields have gridded information of the Earth’s surface properties (e.g.
160 land use, vegetation type, and distribution) and represent surface heterogeneity in the ECLand
161 of the IFS (Kimpson et al., 2023). They are used to compute surface turbulent fluxes (of heat,
162 moisture, and momentum) and skin temperature over different surfaces (vegetation, bare soil,
163 snow, interception, and water) and to calculate an area-weighted average for the grid box for
164 coupling with the atmosphere. To trigger all different parametrization schemes, the ECMWF
165 model uses a set of physiographic fields that do not depend on initial condition of each
166 forecast run or the forecast step. Most fields are constant; surface albedo is specified for 12
167 months to describe the seasonal cycle. Depending on the origin, initial data come at different
168 resolutions and different projections and are then first converted to a regular latitude–
169 longitude grid (EPSG:4326) at ~ 1 km at Equator resolution and secondly to a required grid
170 and resolution. Surface physiographic fields used in this work consist of orographic, land,
171 water, vegetation, soil, albedo fields, see Table 1 for the full list of surface physiographic
172 fields; for more details, see IFS documentation (ECMWF, 2023).

173

174 **2.2.2 ERA5**

175

176 Climate reanalyses combine observations and modelling to provide calculated values of a
 177 range of climactic variables over time. ERA5 is the fifth-generation reanalysis from
 178 ECMWF. It is produced via 4D-Var data assimilation of the IFS cycle 41R2 coupled to a land
 179 surface model (ECLand, (Boussetta et al., 2021)), which includes lake parametrization by
 180 Flake (Mironov and Helmert, n.d.) and an ocean wave model (WAM). The resulting data
 181 product provides hourly values of climatic variables across the atmosphere, land, and ocean
 182 at a resolution of approximately 31 km with 137 vertical sigma levels up to a height of 80 km.
 183 Additionally, ERA5 provides associated uncertainties of the variables at a reduced 63 km
 184 resolution via a 10-member ensemble of data assimilations. In this work, ERA5 hourly
 185 surface fields at ~ 31 km resolution on the cubic octahedral reduced Gaussian grid (i.e.
 186 Tco399) are used. The Gaussian grid's spacing between latitude lines is not regular, but lines
 187 are symmetrical along the Equator; the number of points along each latitude line defines
 188 longitude lines, which start at longitude 0 and are equally spaced along the latitude line. In a
 189 reduced Gaussian grid, the number of points on each latitude line is chosen so that the local
 190 east–west grid length remains approximately constant for all latitudes (here, the Gaussian
 191 grid is N320, where N is the number of latitude lines between a pole and the Equator).

192

193 *Table 1 Input and target features to all emulators from the data sources. The left column*
 194 *shows the observation-derived static physiographic fields, the middle column ERA5 dynamic*
 195 *physiographic and meteorological fields and the rightmost column ECLand generated*
 196 *dynamic prognostic state variables.*

Climate fields	Units	Atmospheric forcing	Units	Prognostic states	Units
Vegetation cover (<i>low, high</i>)		Total precipitation fraction (<i>convective</i> + <i>stratiform</i>)		Soil water volume (<i>Layers</i> <i>1-3</i>)	m^3m^{-3}
Type of vegetation (<i>low, high</i>)		Downward radiation (<i>long,</i> <i>short</i>)	W/m^2	Soil temperature (<i>Layers 1-3</i>)	K
Minimum stomatal resistance (<i>low,</i> <i>high</i>)		Seasonal LAI (<i>high,</i> <i>low</i>)		Snow cover fraction	%

Roughness length (<i>low, high</i>)	Wind speed (v, u)	m/s
Urban cover	Surface pressure	Pa
Lake cover	Skin temperature	K
Lake depth		
Orography (+ <i>std</i> , + m^2/s^{-2} <i>filtered</i>)	Specific humidity	kg/kg
Photosynthesis pathways	Rainfall rate (<i>total</i>)	kg/m ² s
Soil type	Snowfall rate (<i>total</i>)	kg/m ² s
Glacier mask		
Permanent wilting point		
Field capacity		
Cell area		

197

198 2.3 Emulators

199

200 We compare a long-short term memory neural network (LSTM), extreme gradient boosting
 201 regression trees (XGB) and a feedforward neural network (that we here refer to as multilayer
 202 perceptron, MLP). To motivate this setup and pave the way for discussing effects of (hyper-
 203)parameter choices, a short overview of all approaches is given. All analyses were conducted
 204 in Python. XGB was developed in dmlc’s XGBoost python package¹. The MLP and LSTM
 205 were developed in the PyTorch lightning framework for deep learning². Neural networks
 206 were trained with the Adam algorithm for stochastic optimization (Kingma and Ba, 2017).
 207 Model architectures and algorithmic hyperparameters were selected through combined
 208 Bayesian hyperparameter optimization with the Optuna framework (Akiba et al., 2019) and
 209 additional manual tuning. The Bayesian optimization minimizes the neural network
 210 validation accuracy, specified here as mean absolute error (MAE), over a predefined search
 211 space for free hyperparameters with the Tree-structured Parzen Estimator (Ozaki et al., 2022).

¹ <https://xgboost.readthedocs.io/en/stable/python/index.html>

² <https://lightning.ai/docs/pytorch/stable/>

212 The resulting hyperparameter and architecture choices which were used for the different
213 approaches are listed in the Supplementary Material.

214

215 **2.3.1 MLP**

216

217 For creation of the MLP emulator we work with a feed-forward neural network architecture
218 of connected hidden layers with ReLU activations and dropout layers, model components
219 which are given in detail in the Supplementary Material or in (Goodfellow et al., 2016). The
220 MLP was trained with a learning rate scheduler. L2-regularization was added to the training
221 objective via weight decay. Sizes and width of hidden layers as well as hyperparameters were
222 selected together in the hyperparameter optimization procedure. Instead of forecasting
223 absolute prognostic state variables \mathbf{z}_t , the MLP predicts the 6-hourly increment, $\frac{\widehat{dz}}{dt}$. It is
224 trained on a stepwise rollout prediction of future state variables at a pre-defined lead time at
225 given forcing conditions, see details in the section on optimization.

226

227 **2.3.2 LSTM**

228

229 LSTMs are recurrent networks that consider long-term dependencies in time series through
230 gated units with input and forget mechanisms (Hochreiter and Schmidhuber, 1997). In
231 explicitly providing time-varying forcing and state variables, LSTM cell states serve as long-
232 term memory while LSTM hidden states are the cells' output and pass on stepwise short-term
233 representations stepwise. In short notation (Lees et al., 2022), a one-step ahead forward pass
234 followed by a linear transformation can be formulated as

235

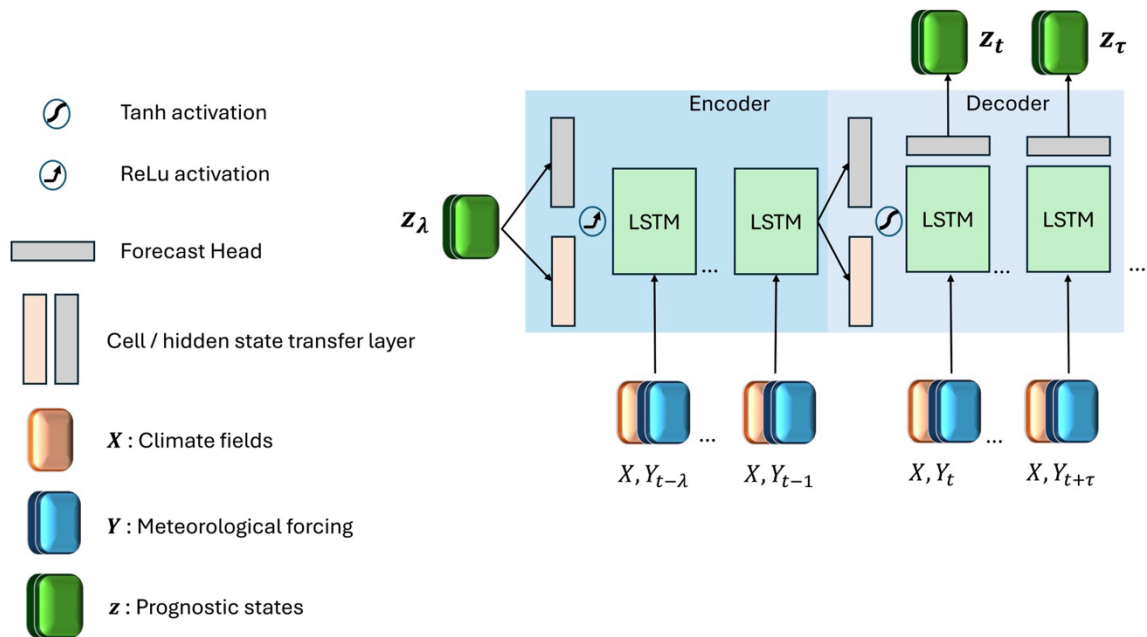
$$\mathbf{h}_t, \mathbf{c}_t = f(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{c}_{t-1}, \boldsymbol{\theta})$$

236

$$\hat{\mathbf{z}}_t = \mathbf{A}\mathbf{h}_t + b$$

237 where \mathbf{h}_{t-1} denotes the hidden state, i.e. output estimates from the previous time step, \mathbf{c}_{t-1}
238 the cell state from the previous time step, and $\boldsymbol{\theta}$ the time-invariant model weights. We stacked
239 multiple LSTM cells to an encoder-decoder model with transfer layers for hidden and cell
240 state initialization and for transfer to the context vector (see figure 1) (Nearing et al., 2024).
241 A lookback l of the previous static and dynamic feature states are passed sequentially to the
242 first LSTM cells in the encoder layer, while the l prognostic state variables \mathbf{z} initialize the
243 hidden state \mathbf{h}_0 after a linear embedding. The output of the first LSTM layer cells become the
244 input to the deeper LSTM layer cells and the last hidden state estimates are the final output

245 from the encoder. Followed by a non-linear transformation with hyperbolic tangent
 246 activation, the hidden cell states are transformed into a weighted context vector \mathbf{s} . Together
 247 with the encoder the cell state (\mathbf{c}_t, \mathbf{s}) initializes the hidden and cell states of the decoder. The
 248 decoder LSTM cells take as input again static and dynamic features sequentially at lead times
 249 $t = 1, \dots, \tau$, but not the prognostic states variables. These are estimated from the sequential
 250 hidden states of the last LSTM layer cells, transformed to target size with a linear forecast
 251 head before prediction. LSTM predicts absolute state variables \mathbf{z}_t while being optimized on
 252 \mathbf{z}_t and $d\hat{\mathbf{z}}_t$ simultaneously, see section on optimization.



253
 254 *Figure 1: LSTM architecture. Blue shaded area indicates the encoder part, where the model*
 255 *is driven by a lookback λ of meteorological forcing and state variables. The light-blue shaded*
 256 *area indicates the decoder part that is initialized from the encoding to unroll LSTM forecasts*
 257 *from the initial time step t up to a flexibly long lead time of τ .*

258 2.3.3 XGB

259
 260 Extreme gradient boosting (XGB) is a regression tree ensemble method that uses an
 261 approximate algorithm for best split finding. It computes first and second order gradient
 262 statistics in the cost function, performing a similar to gradient descent optimization (Chen and
 263 Guestrin, 2016), where each new learner is trained on the residuals of the previous ones.
 264 Regularization and column sampling aim for preventing overfitting internally. XGB is known
 265 to provide a powerful benchmark for time series forecasting and tabular data (Chen and
 266 Guestrin, 2016; Chen et al., 2020; Shwartz-Ziv and Armon, 2021). Like the MLP, it is trained

267 to predict the increment $\widehat{dz}_{t,i}$ of prognostic state variables, but only for a one-step ahead
268 prediction.

269

270 **2.4 Experimental setup**

271

272 We distinguish the experimental analysis into three parts that vary in the usage of the training
273 database: (1) model development, (2) model testing, and (3) global model transfer.

274 The models were developed and for the first time evaluated on a low state resolution

275 (ECMWF’s TCO199 reduced gaussian grid, see section on data sources) and temporal subset

276 from the training data base, i.e. on a bounding box of 7715 grid cells over Europe with time

277 series of six years from 2016 to 2022. For details on the development data base, model

278 selection and model performances, see Supplementary Material S3.

279 The selected models were recreated on a high state resolution (TCO399) continental scale

280 European subset with 10 051 grid cells. Models were trained on five years 2015-2020 with

281 the year 2020 as validation split and evaluated on the year 2021 for the scores we report in

282 the main part. Note that for computation of forecast horizons, the two test years 2021 and

283 2022 were used, see details in section on forecast horizons. With this same data splitting

284 setup, the analysis was repeated in transferring the candidates to the low resolution (TCO199)

285 global data set with a total of 47892 grid cells. The low global resolution on one hand

286 allowed a systematic comparison of the three models, because high resolution training with

287 XGB was prohibited by the required working memory. On the other hand, this extrapolation

288 scenario created an unseen problem for the models that were selected on a continental and

289 high-resolution scale which is reflected in the resulting scores.

290

291 **2.5 Optimization**

292

293 **2.5.1 Loss functions**

294

295 The basis of the loss function \mathcal{L} for the neural network optimization was PyTorch’s

296 SmoothL1Loss³, a robust loss function that combines L1-norm and L2-norm and is less

297 sensitive to outliers than pure L1-norm (Girshick, 2015). Based on a pre-defined threshold

298 parameter β , smooth L1 transitions from L2-norm to L1-norm above the threshold.

³ <https://pytorch.org/docs/stable/generated/torch.nn.SmoothL1Loss.html>

299 SmoothL1Loss \mathcal{L} is defined as

$$300 \quad \mathcal{L}(\hat{z}, z) = 0.5(\hat{z} - z)^2 \frac{1}{\beta} \text{ if } |\hat{z} - z| < \beta \text{ and}$$

$$301 \quad \mathcal{L}(\hat{z}, z) = |\hat{z} - z| - 0.5 \beta \text{ otherwise,}$$

302 here with $\beta = 1$. All models were trained to minimize the incremental loss \mathcal{L}_s that is the
 303 differences between the estimates of the seven prognostic states increments $\widehat{d\mathbf{z}}_t$ and the full
 304 model’s prognostic states increments $d\mathbf{z}_t$ simultaneously as the sum of losses over all states.
 305 We opted for a loss function equally weighted by variables to share inductive biases among
 306 the non-independent prognostic states (Sener and Koltun, 2018). When aggregating over all
 307 training lead times $t = 1, \dots, \tau$, \mathcal{L}_s and grid cells $i = 1, \dots, p$ is

$$308 \quad \mathcal{L}_s(\widehat{d\mathbf{z}}, d\mathbf{z}) = \sum_t^\tau \sum_i^p \mathcal{L}_t(\widehat{d\mathbf{z}}_{t,i}, d\mathbf{z}_{t,i}),$$

309 Whereas when computing a rollout loss \mathcal{L}_r stepwise,

310

$$311 \quad \mathcal{L}_r(\widehat{d\mathbf{z}}, \mathbf{z}) = \frac{1}{\tau} \sum_t^\tau \sum_i^p \mathcal{L}_t(z_{t-1,i} + \widehat{d\mathbf{z}}_{t,i}, z_{t,i})$$

312

313 Prognostic state increments are essentially the first differences from one to the next timestep
 314 that are normalized again by the global standard deviation of the model’s states increments,
 315 s_{dz} before computation of the loss (Keisler, 2022). Due to the forecast models’ structural
 316 differences, loss functions were individually adapted:

317 **MLP** The combined loss function for the MLP is the sum of the incremental loss \mathcal{L}_s and the
 318 rollout loss \mathcal{L}_r . For the rollout loss \mathcal{L}_r , \mathcal{L} was aggregated over grid cells p and accumulated
 319 after an auto-regressive rollout over lead times τ , before being averaged out by division by τ
 320 (Keisler, 2022).

321 **LSTM** The combined loss function for the LSTM is the sum of the incremental loss
 322 \mathcal{L}_s , where the $d\widehat{\mathbf{z}}_t$ were derived from $\widehat{\mathbf{z}}_t$ after the forward pass, and the loss \mathcal{L} computed on
 323 decoder estimates of prognostic states variables, a functionality that leverages the potential of
 324 our LSTM structure.

325 **XGB** Trained only from one to the next time step, i.e. at a lead time of $\tau = 1$, the incremental
 326 loss $\mathcal{L}_s = \mathcal{L}_r$. Without a SmoothL1Loss implementation provided in dmlc’s XGBoost, we
 327 trained XGB with both the Huber-Loss and the default L2-loss. The latter initially providing

328 better results, we chose the default L2-norm as loss function for XGB with the regularization
329 parameter $\lambda = 1$.

330

331 **2.5.1 Normalization**

332 As prognostic target variables are all lower bounded by zero, we tested both z-scoring and
333 max-scoring. The latter yielded no significant improvement; thus we show our results with z-
334 scored target variables. For neural network training but not for fitting XGB, static, dynamic
335 and prognostic state variables were all normalized with z-scoring towards the continental or
336 global mean \bar{z} and unit standard deviation s_z as

$$337 z_{t,n} = \frac{z_{t,n} - \bar{z}}{s_z}.$$

338 Prognostic target state increments were normalized again by the global standard deviation of
339 increments computing the loss (see section 2.5.1) to smooth magnitudes of increments
340 (Keisler, 2022). State variables were back transformed to original scale before evaluation.

341

342 **2.5.3 Spatial and temporal sampling**

343 Sequences were sampled randomly from the training data set, while validation happened
344 sequentially. MLP and XGB were trained on all grid cells simultaneously in both the
345 continental and global setting, while LSTM was trained on the full continental data set but
346 was limited by GPU memory in the global task. We overcame this limitation by randomly
347 subsetting grid cells in the training data into largest possible, equally sized subsets which
348 were then loaded along with the temporal sequences during the batch sampling.

349

350 **2.6 Evaluation**

351

352 Three scores are used for model validation during the model development phase and in
353 validating architecture and hyperparameter selection, being the root mean squared error
354 (*RMSE*), the mean absolute error (*MAE*) and the anomaly correlation coefficient (*ACC*).
355 First, scores were assessed objectively in quantifying forecast accuracy of the emulators
356 against ECLand simulations directly with RMSE and MAE. Doing so, scores were
357 aggregated over lead times, grid cells or both. The total RMSE was computed as

$$358 \text{RMSE} = \sqrt{\frac{\sum_{\tau,p} (z - \hat{z})^2}{n}},$$

359 As the mean absolute error in prognostic state variable prediction over the total of n grid cells
 360 p times lead times τ . Equivalently, MAE was computed as

$$361 \quad \text{MAE} = \frac{\sum_{t,p} |z - \hat{z}|}{n},$$

362 Beyond accuracy, the forecast skill of emulators was assessed using a benchmark model: the
 363 ACC (see below) as index of the long-term naïve climatology c of ECLand, forced by ERA5
 364 (see section 2.2). More specifically, this is the 6-hourly mean of prognostic state variables
 365 over the last 10 years preceding the test year, i.e. the years 2010 to 2020. While climatology
 366 is a hard-to-beat benchmark specifically in long-term forecasting, the persistence is a
 367 benchmark for short-term forecasting (Pappenberger et al., 2015). For verification against
 368 climatology, we compute the anomaly correlation coefficient (ACC) over lead times as

$$369 \quad \text{ACC}(t) = \frac{\overline{(\hat{z} - c)(z - c)}}{\sqrt{\overline{(\hat{z} - c)^2} \overline{(z - c)^2}}}$$

370 at each $t = 1, \dots, \tau$ where the overbar denotes averaging over grid cells $p = i, \dots, n$. This way,
 371 the nominator represents the average spatial covariance of emulator and numerical forecasts
 372 with climatology as expected sample mean. Hence, it indicates the mean squared skill error
 373 towards climatology, and the denominator indicates its variability. The aggregated scores that
 374 are shown in tables 3-5 represent the temporally arithmetic mean of $\text{ACC}(t)$. ACC is bounded
 375 between 1 and -1, and an ACC of 1 indicates perfect representation of forecast error
 376 variability, an ACC of 0.5 indicates a similar forecast error to that of the climatology, an ACC
 377 of 0 indicates that forecast error variability dominates and the forecast has no value and an
 378 ACC approaching -1 indicates that the forecast has been very unreliable (ECMWF, n.d.).
 379 ACC is undefined when the denominator is zero. This is the case either when mean squared
 380 emulator or ECLand anomaly, or both are zero because forecast and climatology perfectly
 381 align, or because they cancel out at summation to the mean.

382

383 **2.6.1 Forecast horizons**

384

385 Forecast horizons of the emulators are defined by the decomposition of the RMSE
 386 (Bengtsson et al., 2008) into the emulator's variability around climatology (i.e. anomaly),
 387 ECLand's variability around climatology and the covariance of both. The horizon is the point
 388 in time at which the forecast error reaches saturation level, that is when the covariance of
 389 emulator and ECLand anomalies approaches zero, as does the ACC.

390 We analysed predictive ability and predictability by computing the ACC for all lead times
391 from 6 hours to approx. one year, i.e. lead times $t = 1, \dots, \tau$, τ being 1350. As this confounds
392 the seasonality with the lead time, we compute these for every starting point of the prediction,
393 requiring two test years (2021 and 2022).

394 Forecast horizons based on the emulators' skill in standardized anomaly towards persistence
395 were equivalently computed but with persistence as a benchmark for shorter time scales, this
396 was only done for three months, from January to March 2021.

397 The analysis was conducted on two exemplary regions in northern and southern Europe that
398 represent very different conditions orography and in prognostic land surface states,
399 specifically in snow cover. For details on the regions and on the horizons computed with
400 standardized anomaly skill, see Appendices A1 and A4 respectively.

401

402 **3 Results**

403

404 The improvement in evaluation runtimes achieved by emulators toward the numerical
405 ECLand were significant. Iterating the forecast over a full test year at 30 km spatial
406 resolution, XGB evaluates in 5.4 minutes, LSTM in 3.09 minutes and MLP in 0.05 minutes
407 (i.e. 3.2 seconds) on average. In contrast, ECLand integration over a full test year on 16
408 CPUs at 30 km spatial resolution takes approximately 240 minutes (i.e. four hours). The slow
409 runtime of the LSTM compared to the MLP emulator is caused by a spatial chunking
410 procedure that was not optimise for this work but could be improved in the future.

411

412 **3.1 Aggregated performances**

413

414 **Europe.** All emulators approximated the numerical LSM with high average total accuracies
415 (all RMSEs < 1.58 and MAEs < 0.84) and confident correlations (all ACC > 0.72) (see table
416 2 and figure 2). The LSTM emulator achieved the best results across all total average scores
417 on the European scale. It decreased the total average MAE by $\sim 25\%$ towards XGB and by
418 $\sim 37\%$ towards the MLP and the total average RMSE by $\sim 42\%$ towards XGB and $\sim 38\%$
419 towards the MLP. In total average ACC, the LSTM scored 20% higher than the MLP and
420 15% than XGB, also being the only emulator that achieved an ACC > 0.9 . While the MLP
421 outperforms XGB in total average RMSE by $\sim 5\%$, XGB scores better than the MLP in MAE
422 by $\sim 27\%$.

423 At variable level, results differentiate into model specific strengths. In soil water volume,
424 XGB outperforms the neural network emulators by up to 60% (m^3m^{-3}) in the first and
425 second layer MAEs towards the LSTM and up to over 40% (m^3m^{-3}) for towards the MLP
426 (see table 3). While the representation of anomalies by specifically the LSTM decreases
427 towards lower soil layers with an ACC of only 0.6214 at the third soil layer, it remains
428 consistently higher for XGB with an ACC still > 0.789 at soil layer three.
429 In soil temperature approximation, LSTM achieves best accuracies at higher soil levels with
430 up to 7% (K) improvement in MAE towards XGB and ACCs > 0.92 , but XGB outperforms
431 LSTM at the third soil level with a close to 50% (K) improvement (see table 4). The MLP
432 doesn't stand out by high scores on the continental scale. However, in terms of accuracy we
433 found an inverse ranking in the model development procedure during which LSTM outscored
434 XGB in soil water volume but struggled with soil temperature approximations, for the
435 interested reader we refer to the supplementary information.
436 In snow cover approximation, the LSTM emulator enhances accuracies by over $\sim 50\%$ in
437 MAE towards both the XGB and the MLP emulator and scores highest in anomaly
438 representation with an ACC of ~ 0.87 compared to an ACC of ~ 0.66 for the MLP and only
439 ~ 0.74 for the XGB (see table 5).
440 **Globe.** Score ranking on the global scale varies strongly from the continental scale (see table
441 2). In total average accuracies, the MLP outperforms XGB by over 30% and LSTM by up
442 $\sim 25\%$ in RMSE and improves MAE more than 15% towards both. In anomaly correlation
443 however it scores last, whereas XGB achieves the highest total average of over 0.75.
444 Consistent with scores on the continental scale is XGBs high performance in soil temperature
445 (see table 3). It significantly outperforms the LSTM by $\sim 60\%$ (K) in RMSE and nearly up to
446 75% (K) in MAE in all layers and the MLP by up to 50% (K) in MAE at the top layer.
447 Anomaly persistence for all models degrade visibly towards the lower soil layers, while that
448 of the LSTM most relative to MLP and XGB. Like on the continental scale, XGB also
449 outperforms the other candidates in soil temperature forecasts in all but the medium layer,
450 where the MLP gets higher scores in MAE and RMSE but not in ACC (see table 4). LSTM
451 doesn't stand out with any scores on the global scale.

452

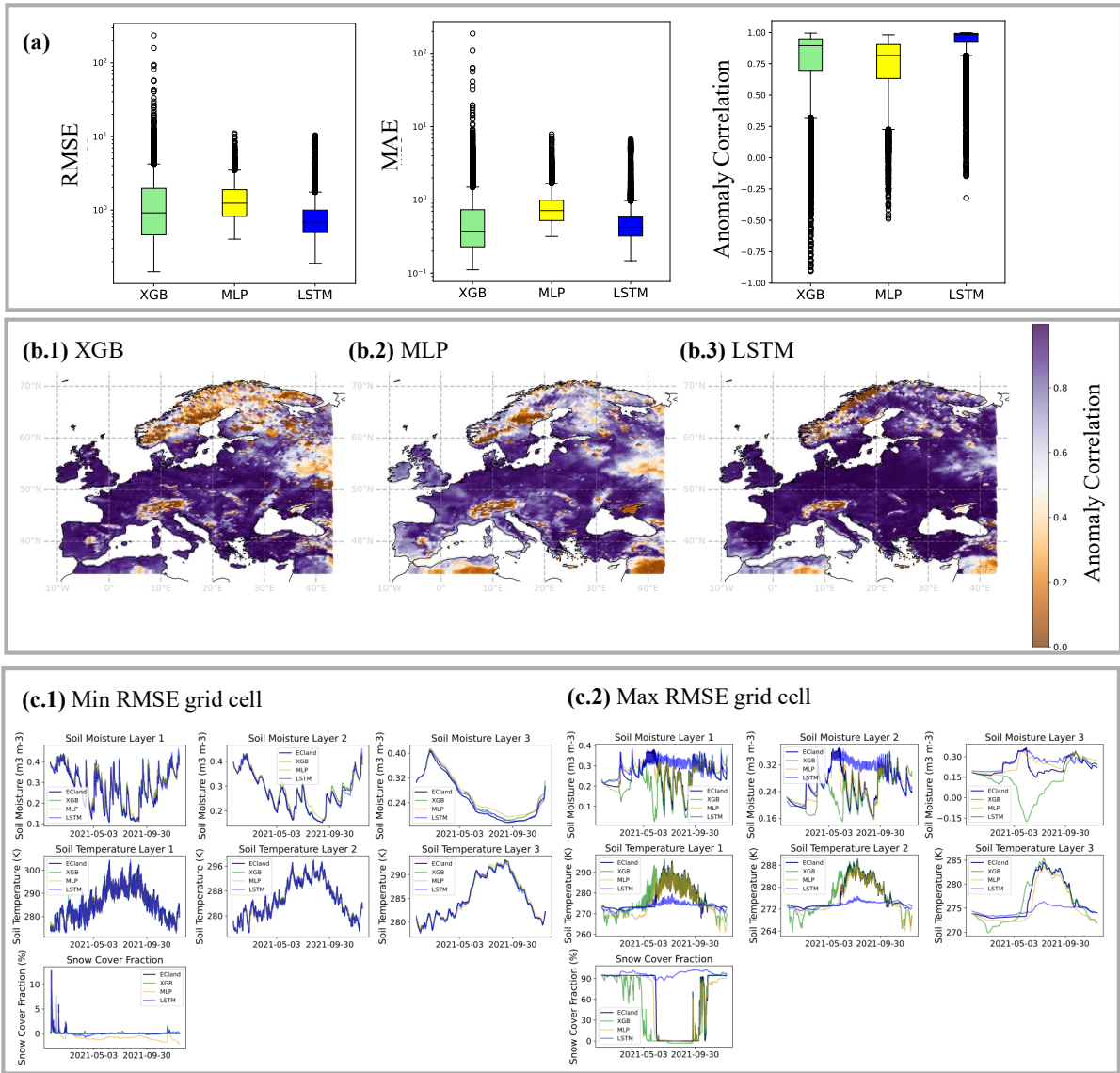
453 3.2 Spatial and temporal performances

454

455 **Europe.** When summarizing temporally aggregated scores as boxplots to a total distribution
456 over space (see figure 2, A), the long tails of XGB scores become visible, whereas the MLP

457 indicates most robustness. This is reflected in the geographic distribution of scores at the
458 example of ACC (see figure 2, bottom), where the area of low anomaly correlation is largest
459 for XGB, ranging over nearly all northern Scandinavia, while MLP and LSTM have smaller
460 and more segregated areas of clearly low anomaly correlation. The LSTM shows a
461 homogenously high ACCs over most of central Europe but the Alps, while also seems to be
462 challenged in areas of relative to the central Europe extreme weather conditions at the
463 Norwegian and Spanish coasts.

464 **Globe.** Like the results from the continental analysis, we find again long upper tails of
465 outliers for XGB in total spatial distribution of accuracies, both in RMSE and MAE and only
466 few outliers for MLP and LSTM. The anomaly correlation distribution changed towards
467 longer lower tails for MLP and LSTM and a shorter lower tail for XGB. We should, however,
468 take the results of total average ACC with care as it remains largely undefined in regions
469 without much noise in snow cover or soil water volume and globally represents mainly
470 patterns of soil temperature.



471

472 *Figure 2: a: Total aggregated distributions of (log) scores averaged over lead times, i.e.*
 473 *displaying the variation among grid cells. b: The distribution of the anomaly correlation in*
 474 *space on the European subset (b.1: XGB, b.2: MLP, b.3: LSTM). c: Model forecasts over test*
 475 *year 2021 for grid cell with minimum and maximum RMSE values (LSTM).*

476

477 *Table 2: Emulator total average scores (unitless), aggregated over variables, time and space*
 478 *from the European and Global model testing.*

Variable	Model	RMSE		MAE		ACC	
		Europe	Globe	Europe	Globe	Europe	Globe
All variables	XGB	1.575	2.611	0.695	1.601	0.765	0.755
	MLP	1.486	1.699	0.832	1.189	0.728	0.569
	LSTM	0.918	2.252	0.526	1.787	0.925	0.647

479 *Table 3: Emulator average scores (RMSE, MAE in m^3m^{-3}) on soil water volume forecasts*
 480 *for the European subset, aggregated over space and time from the European and Global*
 481 *model testing.*

Variable	Layer	Model	RMSE		MAE		ACC	
			Europe	Globe	Europe	Globe	Europe	Globe
Soil water volume	1	XGB	0.013	0.015	0.01	0.01	0.908	0.92
		MLP	0.019	0.029	0.015	0.023	0.856	0.791
		LSTM	0.029	0.048	0.023	0.04	0.847	0.729
	2	XGB	0.011	0.012	0.008	0.009	0.901	0.884
		MLP	0.019	0.023	0.014	0.018	0.789	0.77
		LSTM	0.029	0.05	0.023	0.042	0.79	0.617
	3	XGB	0.015	0.014	0.011	0.01	0.789	0.777
		MLP	0.02	0.02	0.017	0.016	0.576	0.667
		LSTM	0.033	0.051	0.027	0.043	0.621	0.475

482

483 *Table 4: Emulators' average scores (RMSE, MAE in K) on soil temperature forecasts for the*
 484 *European subset, aggregated over space and time.*

Variable	Layer	Model	RMSE		MAE		ACC	
			Europe	Globe	Europe	Globe	Europe	Globe
Soil temperature	1	XGB	1.154	4.539	0.744	3.278	0.806	0.769
		MLP	1.628	2.606	1.188	2.072	0.674	0.581
		LSTM	0.931	3.152	0.682	2.626	0.938	0.735
	2	XGB	0.901	2.501	0.51	1.772	0.812	0.797
		MLP	1.134	1.851	0.784	1.452	0.718	0.606
		LSTM	0.734	2.87	0.541	2.4	0.928	0.699
	3	XGB	0.714	1.287	0.482	0.933	0.722	0.711
		MLP	1.128	1.375	0.821	1.071	0.416	0.514
		LSTM	1.141	3.466	0.918	3.002	0.598	0.406

485

486 *Table 5: Emulators' average scores (RMSE, MAE in %) on snow cover forecasts for the*
 487 *European subset, aggregated over space and time.*

Variable	Layer	Model	RMSE		MAE		ACC	
			Europe	Globe	Europe	Globe	Europe	Globe

Snow cover	top	XGB	8.219	9.906	3.099	5.196	0.746	0.707
		MLP	6.449	5.995	2.986	3.671	0.66	0.618
		LSTM	3.526	6.127	1.47	4.357	0.877	0.698

488

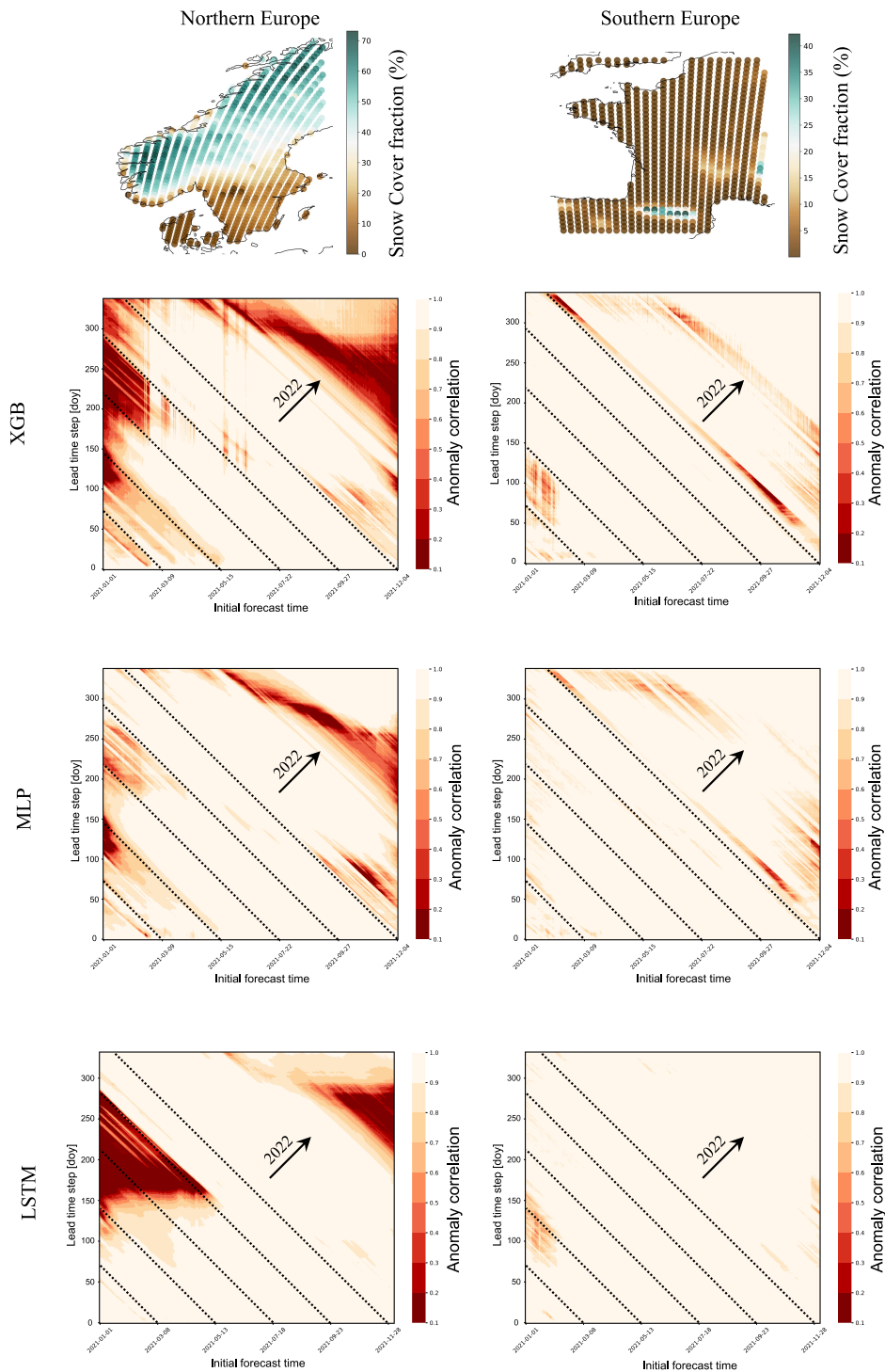
489 3.3 Forecast horizons

490 Forecast horizons were computed for two European regions, of which the northern one
491 represents the area of lowest emulators' skill (see figure 2, B.1-3) and the southern one an
492 area of stronger emulators' skill. Being strongly correlated with soil water volume, these two
493 regions differ specifically in their average snow cover fraction (see figure 3). The displayed
494 horizons were computed over all prognostic state variables simultaneously, while their
495 interpretation is related to horizons computed for prognostic state variables separately, for the
496 figures of which we refer to the Supplementary Material.

497 In the North, predictive skill depended on an interaction of how far ahead a prediction was
498 made (the lead time) and the day of year to which the prediction was made. In the best case,
499 the LSTM, summer predictions were poor (light patches in figure 3 heat maps), but only
500 when initialised in winter. Or, in other words, one can make good predictions starting in
501 winter, but not to summer. Vertical structures indicate a systematic model error that appears at
502 specific initialisation times and that is independent of prediction date, for example in XGB
503 forecasts that are initialized in May (see figure 3, northern region). Diagonal light structures
504 in the heat maps indicate a temporally consistent error and can be interpreted as physical
505 limits of system predictability, where the different initial forecast time doesn't affect model
506 scores.

507 All models show stronger limits in predictability and predictive ability in the northern
508 European region (see figure 3, left column). MLP and XGB struggled with representing
509 seasonal variation towards climatology at long lead times, while LSTM is strongly limited by
510 a systematic error in certain regions. Initializing the forecast the 1 January 2021, MLP drops
511 below an ACC of 80% repeatedly from initialization on and then to an ACC below 10% at the
512 beginning of May. LSTMs performance is more robust in the beginning of the year but
513 depletes strongly later to less than 10% ACC in mid-May. On the one hand, this represents
514 two different characteristics of model errors: MLP forecasts for snow cover fraction are less
515 than zero for some grid cells while LSTM forecasts for snow cover fraction remain falsely at
516 very high levels for some grid cells, not predicting the snowmelt in May (see Supplementary
517 Material, S4.1). On the other hand, this represents a characteristic error due to change in

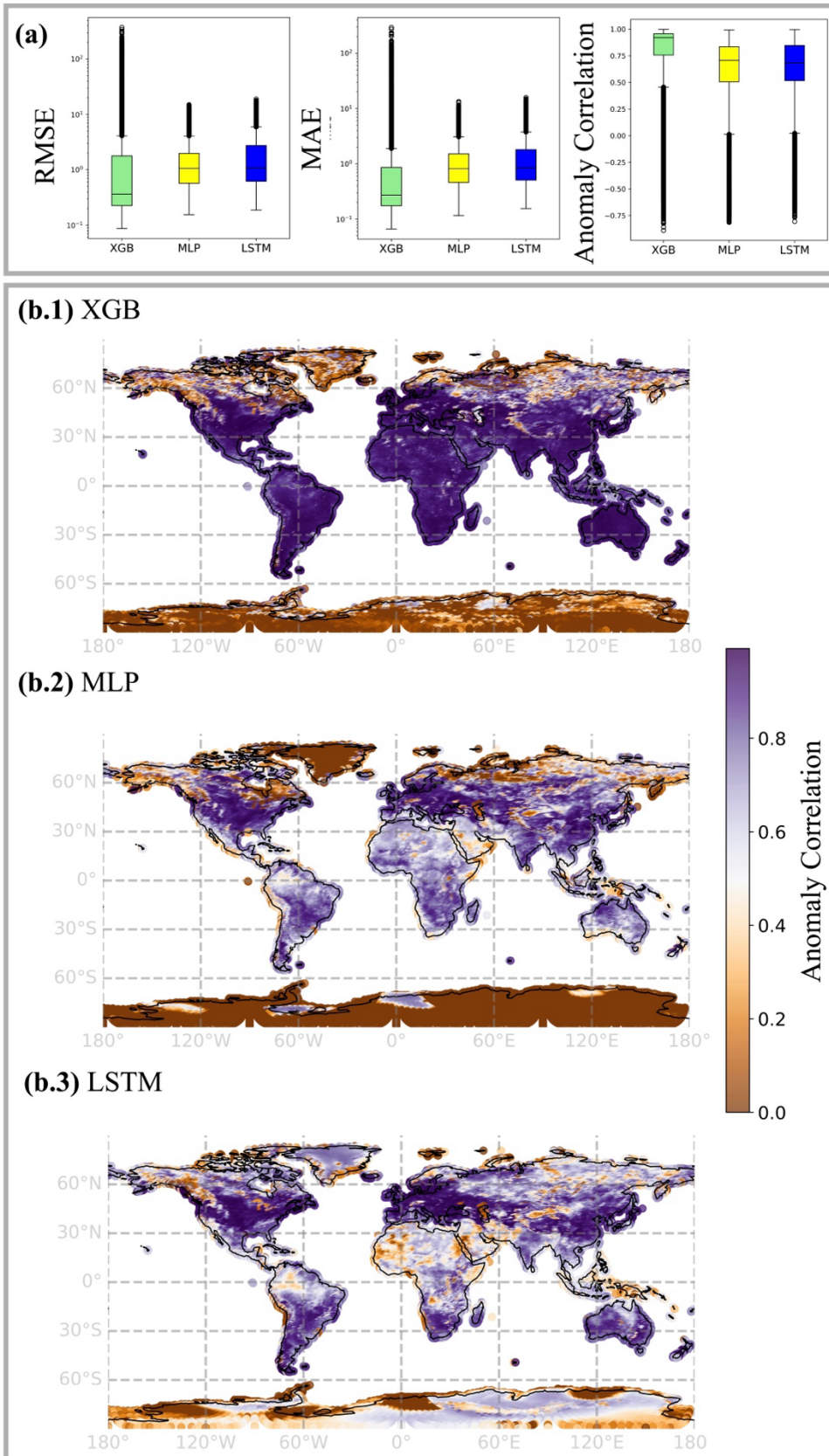
518 seasonality: the snowmelt in this region in May happens abruptly and all emulators
 519 repeatedly over- or underpredict the exact date.
 520



521
 522 *Figure 3: Top row: European subregions for computations of forecast skill horizons and their*
 523 *yearly average snow cover fraction (%), predicted by ECLand. Rows 2-4: Emulator forecast*
 524 *skill horizons in the subregions, aggregated over prognostic state variables, computed with*
 525 *the anomaly correlation coefficient (ACC) at 6-hourly lead times (y-axis) over approx. one*

526 *year, displayed as a function of the initial forecast time (x-axis). The horizon is the time at*
527 *which the forecast has no value at all, i.e. when ACC is 0 (or below 10%). The diagonal*
528 *dashed lines indicate the day of the test year 2021 as labelled on the x-axis, the arrows*
529 *indicate where forecasts reach the second test year 2022.*

530



531

532 *Figure 4: a) Total average scores, representing spatial variation among grid cells. B) Total*
 533 *average ACC in space. Note that ACC remained undefined for regions of low signal in snow*
 534 *cover and soil water volume, see Supplementary Material.*

535 **4 Discussion**

536

537 In the comparative analysis of emulation approaches for land surface forecasting, three
538 primary models—LSTM (Long Short-Term Memory networks), MLP (Multi-Layer
539 Perceptrons), and XGB (Extreme Gradient Boosting)—have been evaluated to understand
540 their effectiveness across different operational scenarios. Evaluating emulators over the test
541 period yielded a significant runtime improvement toward the numerical model for all
542 approaches (see section 3). While all models achieved high predictive scores, they differ in
543 their demand of computational resources (Cui et al., 2021) and each one offers unique
544 advantages and faces distinct challenges, impacting their suitability for various forecasting
545 tasks. In this work we present the first steps towards enabling quick offline experimentation
546 on the land surface with ECMWF’s land surface scheme ECLand and towards decreasing
547 computational demands in, i.e. coupled data assimilation.

548

549 **4.1 Approximation of prognostic land surface states**

550

551 The total evaluation scores of our emulators indicate good agreement with ECLand
552 simulations. Among the seven individual prognostic land surface states, emulators achieve
553 notably different scores and in the transfer from the high-resolution continental to the low-
554 resolution global scale, their performance ranking change. On average, neural network
555 performances degrade towards the deeper soil layers, while XGB scores remain relatively
556 stable. Also, the neural networks scores drop in the extrapolation from continental to global
557 scale, while XGB scores also for this task remain constantly high.

558 In a way, these findings are not surprising. It is known that neural networks are highly
559 sensitive to selection bias (Grinsztajn et al., 2022) and tuning of hyper-parameters
560 (Bouthillier et al., 2021), suboptimal choices of which may destabilise variance in predictive
561 skill. Previous and systematic comparisons of XGB and deep neural networks have
562 demonstrated that neural networks can hardly be transferred to new data sets without
563 performance loss (Shwartz-Ziv and Armon, 2021). On tabular data, XGB still outperforms
564 neural networks in most cases (Grinsztajn et al., 2022), unless these models are strongly
565 regularized (Kadra et al., 2021). The disadvantage of neural networks might lay in the
566 rotational invariance of MLP-like architectures, due to which information about the data
567 orientation gets lost, as well as in their instability regarding uninformative input features
568 (Grinsztajn et al., 2022).

569 Inversely to expectations and preceding experiments, on the European data set relative to the
570 two other models the LSTM scored better in the upper layer soil temperatures than in
571 forecasting soil water volume and decreased in scores towards lower layers with slower
572 processes. For training on observations, the decreasing LSTM predictive accuracy for soil
573 moisture with lead time is discussed (Datta and Faroughi, 2023), but reasons arising from the
574 engineering side remain unclear. In an exemplary case of a single-objective, deterministic
575 streamflow forecast, a decrease in recurrent neural network performance has been related
576 with an increasing coefficient of variation (Guo et al., 2021). In our European subregions, the
577 signal-to-noise ratio of the prognostic state variables (computed as the averaged ratio of mean
578 and standard deviation) is up to ten times higher in soil temperature than in soil water volume
579 states (see Supplementary Material, S2.1). While a small signal of the latter may induce
580 instability in scores, it does not explain the decreasing performance towards deeper soil layers
581 with slow processes, where we expected an advantage of the long-term memory.
582 Stein's paradox tells us that joint optimization may lead to better results if the target is multi-
583 objective, but not if we are interested in single targets (James and Stein, 1992; Sener and
584 Koltun, 2018). While from a process perspective multi-objective scores are less meaningful
585 than single ones, this is what we opted for due to efficiency. The unweighted linear loss
586 combination might be suboptimal in finding effective parameters across all prognostic state
587 variables (Chen et al., 2017; Sener and Koltun, 2018), yet being strongly correlated, we
588 deemed their manual weighting inappropriate. An alternative to this provides adaptive loss
589 weighting with gradient normalisation (Chen et al., 2017).

590

591 **4.2 Evaluation in time and space**

592

593 We used aggregated MAE and RMSE accuracies as a first assessment tool to conduct model
594 comparison, but score aggregation hides model specific spatio-temporal residual patterns.
595 Further, both scores are variance dependent, favouring low variability in model forecasts
596 even though this may not be representative of the system dynamic (Thorpe et al., 2013).
597 Assessing the forecast skill over time as the relative proximity to a subjectively chosen
598 benchmark helps disentangling areas of strengths and weaknesses in forecasting with the
599 emulators (Pappenberger et al., 2015). The naïve 6-hourly climatology as benchmark
600 highlights periods where emulators long-range forecasts on the test year are externally limited
601 by seasonality, i.e. system predictability, and where they are internally limited by model error,
602 i.e. the model's predictive ability. Applying this strategy in two exemplary European

603 subregions showed that all emulators struggle most in forecasting the period from late
604 summer to autumn, unless they are initialized in summer (see figure 3). Because forecast
605 quality is most strongly limited by snow cover (see Supplementary Material, A4.1), we
606 interpret this as the unpredictable start of snow fall in autumn. External predictability
607 limitations seem to affect the LSTM overall less than the two other models, and specifically
608 XGB drifts at long lead times.

609 From a geographical perspective inferred from the continental scale, emulators struggle in
610 forecasting prognostic state variables in regions with complicated orography and strong
611 environmental gradients. XGB scores vary seemingly random in space, while neural
612 networks scores exhibit spatial autocorrelation. A meaningful inference about this, however,
613 can only be conducted in assessing model sensitivities to physiographic and meteorological
614 fields through gradients and partial dependencies. While the goal of this work is to introduce
615 our approach to emulator development, this can be investigated in future analyses.

616

617 **4.3 Emulation with memory mechanisms**

618

619 Without much tuning, XGB challenges both LSTM and MLP for nearly all variables (see
620 tables 2-4). In training on observations for daily short-term and real-time rainfall-runoff
621 prediction, XGB and LightXGB were shown before to equally performed as, or outperformed
622 LSTMs (Chen et al., 2020; Cui et al., 2021). Nevertheless, models with memory mechanism
623 such as the encoder-decoder LSTM remain a promising approach for land surface forecasting
624 regarding their differentiability (Hatfield et al., 2021), their flexible extension of lead times,
625 for exploring the effect of long-term dependencies or for inference from the context vector
626 that may help identifying the process relevant climate fields (Lees et al., 2022).

627 The LSTM architecture assumes that the model is well defined in that the context vector
628 perfectly informs the hidden decoder states. If that assumption is violated, potential strategies
629 are to create a skip-connection between context vector and forecast head, or to consider input
630 of time-lagged variables or self-attention mechanisms (Chen et al., 2020). With attention, the
631 context vector becomes a weighted sum of alignments that relates neighbouring positions of a
632 sequence, a feature that could be leveraged for forecasting quick processes such as snow
633 cover or top-level soil water volume.

634 Comparing average predictive accuracies across different training lead times indicates that
635 training at longer lead times may enhance short-term accuracy of the LSTM at the cost of
636 training runtime (see Supplementary Material, S2). A superficial exploration of encoder

637 length indicates no visible improvement on target accuracies if not a positive tendency
638 towards shorter sequences. This needs an extended analysis for understanding, yet without a
639 significant improvement by increased sequence length, GRU cells might provide a simplified
640 and less parameterized alternative to LSTM cells. They were found to perform equally well
641 on streamflow forecast performance before, while reaching higher operational speed (Guo et
642 al., 2021).

643

644 **4.4 Emulators in application**

645

646 LSTM networks with a decoder structure are valued for their flexible and fast lead time
647 evaluation, which is crucial in applications where forecast intervals are not consistent. The
648 structure of LSTM is well-suited for handling sequential data, allowing it to perform
649 effectively over different temporal scales (Hochreiter and Schmidhuber, 1997). They provide
650 access to gradients, which facilitates inference, optimization and usage for coupled data
651 assimilation (Hatfield et al., 2021). Nevertheless, the complexity of LSTMs introduces
652 disadvantages: Despite their high evaluation speed and accuracy under certain conditions,
653 they require significant computational resources and long training times. They are also highly
654 sensitive to hyperparameters, making them challenging to tune and slow to train, especially
655 with large datasets.

656 MLP models stand out for their implementation, training and evaluation speed with yet
657 rewarding accuracy, making them a favourable choice for scenarios that require rapid model
658 deployment. They are tractable and easy to handle, with a straightforward setup that is less
659 demanding computationally than more complex models. MLPs also allow for access to
660 gradients, aiding in incremental improvements during training and quick inference (Hatfield
661 et al., 2021). Despite these advantages, MLPs face challenges with memory scaling during
662 training at fixed lead times, which can hinder their applicability in large-scale or high-
663 resolution forecasting tasks.

664 XGB models are highly regarded for their robust performance with minimal tuning,
665 achieving high accuracy not only in sample applications, but also in transfer to unseen
666 problems (Grinsztajn et al., 2022; Shwartz-Ziv and Armon, 2021). Their simplicity makes
667 them easy to handle, even for users with limited technical expertise in machine learning.
668 However, the slow evaluation speed of XGB becomes apparent as dataset complexity and
669 size increase. Although generally more interpretable than deep machine learning tools, XGB

670 is not differentiable, limiting its application in coupled data assimilation (Hatfield et al.,
671 2021) even though research on differentiable trees is ongoing (Popov et al., 2019).

672

673 **4.5 Experimentation with Emulators**

674

675 In the IFS, the land surface is coupled to the atmosphere via skin temperature (ECMWF,
676 2023), the predictability of which is known to be influenced by specifically by soil moisture
677 (Dunkl et al., 2021). This is the interface with the numerical model where a robust surrogate
678 could act online to improve forward (i.e. parametrization (Brenowitz et al., 2020)) or
679 backward (i.e. data assimilation (Hatfield et al., 2021)) procedures, and it motivates the
680 experiment from the perspective of hybrid forecasting models (Irrgang et al., 2021; Slater et
681 al., 2023). However, because an offline training ignores the interaction with the atmospheric
682 model, emulator scores will not directly translate to the coupled performance and of course
683 additional experiments would be necessary (Brenowitz et al., 2020). As the current stand-
684 alone models, emulators provide a pre-trained model-suite (Gelbrecht et al., 2023) and can be
685 used for experimentation on the land surface. The computation of forecast horizons is an
686 example for such an experiment, seen as a step toward a predictability analysis of land
687 surface processes. Full predictability analyses are commonly conducted with model
688 ensembles (Guo et al., 2011; Shukla, 1981), the simulation of which can quicker be
689 done with emulators than with the numerical model (see evaluation runtimes, section
690 3).

691 We want to stress at this point that to avoid misleading statements, evaluation of the
692 emulators on observations is required. In the context of surrogate models, two inherent
693 sources of uncertainty are specifically relevant: First, the structural uncertainty by
694 statistical approximation of the numerical model and second, the uncertainty arising by
695 parameterization with synthetic (computer model generated) data (Brenowitz et al.,
696 2020; Gu et al., 2017). Both sources can cause instabilities in surrogate models that
697 could translate when coupled with the IFS (Beucler et al., 2021), but that also should be
698 quantified when drawing conclusions from the stand-alone models outside of the
699 synthetic domain. Consequently, a reliable surrogate model for online or offline
700 experimentation requires validation, and enforcing additional constraints may be
701 advantageous for physical consistency (Beucler et al., 2021).

702

703 **5 Conclusion**

704

705 To conclude, the choice between LSTM, MLP, and XGB models for land surface forecasting
706 depends largely on the specific requirements of the application, including the need for speed,
707 accuracy, and ease of use. Each model's computational demands, flexibility, and operational
708 overhead must be carefully considered to optimize performance and applicability in diverse
709 forecasting environments. When it comes to accuracy, combined model ensembles of XGB
710 and neural networks have been shown to yield the best results (Shwartz-Ziv and Armon,
711 2021), but accuracy alone will not determine a single best approach (Bouthillier et al., 2021).
712 Our comparative assessment underscores the importance of selecting the appropriate
713 emulation approach based on a clear understanding of each model's strengths and limitations
714 in relation to the forecasting tasks at hand. By developing the emulators for ECMWF's
715 numerical land surface scheme ECLand, we path the way towards a physics-informed ML-
716 based land surface model that on the long run can be parametrized with observations. We also
717 provide a pretrained model suite to improve land surface forecasts and future land reanalyses.

718

719 **Code and data availability**

720 Code for this analysis is published [on OSF \(DOI: 10.17605/OSF.IO/8567D\)](https://doi.org/10.17605/OSF.IO/8567D) or at
721 <https://github.com/MWesselkamp/land-surface-emulation>. Training data is published at
722 [10.21957/n17n-6a68](https://doi.org/10.21957/n17n-6a68) (Tco199) and [10.21957/pcf3-ah06](https://doi.org/10.21957/pcf3-ah06) (Tco399).

723 **Author contribution**

724 MW, MCha, EP, FP and GB conceived the study. MW and EP conducted the analysis. MW,
725 MCha, MK, EP discussed and took technical decisions. SB advised on process decisions.
726 MW, MCho and FP wrote the manuscript. MW, MCha, EP, MCho, SB, MK, CFD, FP
727 reviewed the analysis and/or manuscript.

728 **Competing interest**

729 The authors declare that they have no conflict of interest.

730 **Acknowledgements**

731 This work profited from discussion with Linus Magnusson, Patricia de Rosnay, Sina R. K.
732 Farhadi and Karan Ruparell and many more. MW thankfully acknowledges ECMWF for
733 providing two research visit stipendiatees over the course of the collaboration. EP was funded
734 by the CERISE project (grant agreement No101082139) funded by the European Union.
735 Views and opinions expressed are however those of the authors only and do not necessarily

736 reflect those of the European Union or the Commission. ChatGPT version 4.0 was used for
737 coding support.

738

739 **References**

740

741 Akiba, T., Sano, S., Yanase, T., Ohta, T., and Koyama, M.: Optuna: A Next-generation
742 Hyperparameter Optimization Framework, in: Proceedings of the 25th ACM SIGKDD
743 International Conference on Knowledge Discovery & Data Mining, KDD '19: The 25th
744 ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Anchorage AK
745 USA, 2623–2631, <https://doi.org/10.1145/3292500.3330701>, 2019.

746 Baker, E., Harper, A. B., Williamson, D., and Challenor, P.: Emulation of high-resolution
747 land surface models using sparse Gaussian processes with application to JULES,
748 *Geosci. Model Dev.*, 15, 1913–1929, <https://doi.org/10.5194/gmd-15-1913-2022>, 2022.

749 Balsamo, G., Boussetta, S., Dutra, E., Beljaars, A., and Viterbo, P.: Evolution of land-
750 surface processes in the IFS, 2011.

751 Bassi, A., Höge, M., Mira, A., Fenicia, F., and Albert, C.: Learning Landscape Features
752 from Streamflow with Autoencoders, <https://doi.org/10.5194/hess-2024-47>, 20
753 February 2024.

754 Bengtsson, L. K., Magnusson, L., and Källén, E.: Independent Estimations of the
755 Asymptotic Variability in an Ensemble Forecast System, *Mon. Weather Rev.*, 136, 4105–
756 4112, <https://doi.org/10.1175/2008MWR2526.1>, 2008.

757 Beucler, T., Pritchard, M., Rasp, S., Ott, J., Baldi, P., and Gentine, P.: Enforcing Analytic
758 Constraints in Neural Networks Emulating Physical Systems, *Phys. Rev. Lett.*, 126,
759 098302, <https://doi.org/10.1103/PhysRevLett.126.098302>, 2021.

760 Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X., and Tian, Q.: Pangu-Weather: A 3D High-
761 Resolution Model for Fast and Accurate Global Weather Forecast,
762 <https://doi.org/10.48550/ARXIV.2211.02556>, 2022.

763 Boussetta, S., Balsamo, G., Arduini, G., Dutra, E., McNorton, J., Choulga, M., Agustí-
764 Panareda, A., Beljaars, A., Wedi, N., Munõz-Sabater, J., De Rosnay, P., Sandu, I.,
765 Hadade, I., Carver, G., Mazzetti, C., Prudhomme, C., Yamazaki, D., and Zsoter, E.:
766 ECLand: the ECMWF land surface modelling system, *Atmosphere*, 12, 723,
767 <https://doi.org/10.3390/atmos12060723>, 2021.

768 Bouthillier, X., Delaunay, P., Bronzi, M., Trofimov, A., Nichyporuk, B., Szeto, J., Sepah, N.,
769 Raff, E., Madan, K., Voleti, V., Kahou, S. E., Michalski, V., Serdyuk, D., Arbel, T., Pal, C.,
770 Varoquaux, G., and Vincent, P.: Accounting for Variance in Machine Learning
771 Benchmarks, <http://arxiv.org/abs/2103.03098>, 1 March 2021.

- 772 Brenowitz, N. D., Henn, B., McGibbon, J., Clark, S. K., Kwa, A., Perkins, W. A., Watt-
773 Meyer, O., and Bretherton, C. S.: Machine Learning Climate Model Dynamics: Offline
774 versus Online Performance, <https://doi.org/10.48550/ARXIV.2011.03081>, 2020.
- 775 Chantry, M., Hatfield, S., Duben, P., Polichtchouk, I., and Palmer, T.: Machine learning
776 emulation of gravity wave drag in numerical weather forecasting,
777 <https://doi.org/10.48550/ARXIV.2101.08195>, 2021.
- 778 Chen, T. and Guestrin, C.: XGBoost: A Scalable Tree Boosting System,
779 <https://doi.org/10.48550/ARXIV.1603.02754>, 2016.
- 780 Chen, X., Huang, J., Han, Z., Gao, H., Liu, M., Li, Z., Liu, X., Li, Q., Qi, H., and Huang, Y.:
781 The importance of short lag-time in the runoff forecasting model based on long short-
782 term memory, *J. Hydrol.*, 589, 125359, <https://doi.org/10.1016/j.jhydrol.2020.125359>,
783 2020.
- 784 Chen, Z., Badrinarayanan, V., Lee, C.-Y., and Rabinovich, A.: GradNorm: Gradient
785 Normalization for Adaptive Loss Balancing in Deep Multitask Networks,
786 <https://doi.org/10.48550/ARXIV.1711.02257>, 2017.
- 787 Choulga, M., Kourzeneva, E., Balsamo, G., Boussetta, S., and Wedi, N.: Upgraded global
788 mapping information for earth system modelling: an application to surface water depth
789 at the ECMWF, *Hydrol. Earth Syst. Sci.*, 23, 4051–4076, <https://doi.org/10.5194/hess-23-4051-2019>, 2019.
- 791 Cui, Z., Qing, X., Chai, H., Yang, S., Zhu, Y., and Wang, F.: Real-time rainfall-runoff
792 prediction using light gradient boosting machine coupled with singular spectrum
793 analysis, *J. Hydrol.*, 603, 127124, <https://doi.org/10.1016/j.jhydrol.2021.127124>, 2021.
- 794 Datta, P. and Faroughi, S. A.: A multihead LSTM technique for prognostic prediction of
795 soil moisture, *Geoderma*, 433, 116452,
796 <https://doi.org/10.1016/j.geoderma.2023.116452>, 2023.
- 797 De Rosnay, P., Balsamo, G., Albergel, C., Muñoz-Sabater, J., and Isaksen, L.:
798 Initialisation of Land Surface Variables for Numerical Weather Prediction, *Surv.*
799 *Geophys.*, 35, 607–621, <https://doi.org/10.1007/s10712-012-9207-x>, 2014.
- 800 Dunkl, I., Spring, A., Friedlingstein, P., and Brovkin, V.: Process-based analysis of
801 terrestrial carbon flux predictability, *Earth Syst. Dyn.*, 12, 1413–1426,
802 <https://doi.org/10.5194/esd-12-1413-2021>, 2021.
- 803 ECMWF: IFS Documentation CY43R3 - Part IV: Physical processes,
804 <https://doi.org/10.21957/EFYK72KL>, 2017.
- 805 ECMWF: IFS Documentation CY48R1 - Part IV: Physical Processes,
806 <https://doi.org/10.21957/02054F0FBB>, 2023.
- 807 ECMWF: Forecast User Guide, *Anom. Correl. Coeff.*, n.d.

808 Fer, I., Kelly, R., Moorcroft, P. R., Richardson, A. D., Cowdery, E. M., and Dietze, M. C.:
809 Linking big models to big data: efficient ecosystem model calibration through Bayesian
810 model emulation, *Biogeosciences*, 15, 5801–5830, [https://doi.org/10.5194/bg-15-5801-](https://doi.org/10.5194/bg-15-5801-2018)
811 2018, 2018.

812 Gelbrecht, M., White, A., Bathiany, S., and Boers, N.: Differentiable programming for
813 Earth system modeling, *Geosci. Model Dev.*, 16, 3123–3135,
814 <https://doi.org/10.5194/gmd-16-3123-2023>, 2023.

815 Girshick, R.: Fast R-CNN, <https://doi.org/10.48550/ARXIV.1504.08083>, 2015.

816 Goodfellow, I., Bengio, Y., and Courville, A.: Deep learning, The MIT Press, Cambridge,
817 Massachusetts, 775 pp., 2016.

818 Grinsztajn, L., Oyallon, E., and Varoquaux, G.: Why do tree-based models still
819 outperform deep learning on tabular data?, <https://doi.org/10.48550/ARXIV.2207.08815>,
820 2022.

821 Gu, M., Wang, X., and Berger, J. O.: Robust Gaussian Stochastic Process Emulation,
822 <https://doi.org/10.48550/ARXIV.1708.04738>, 2017.

823 Guo, Y., Yu, X., Xu, Y.-P., Chen, H., Gu, H., and Xie, J.: AI-based techniques for multi-step
824 streamflow forecasts: application for multi-objective reservoir operation optimization
825 and performance assessment, *Hydrol Earth Syst Sci*, 2021.

826 Guo, Z., Dirmeyer, P. A., and DelSole, T.: Land surface impacts on subseasonal and
827 seasonal predictability: LAND IMPACTS SUBSEASONAL PREDICTABILITY, *Geophys. Res.*
828 *Let.*, 38, n/a-n/a, <https://doi.org/10.1029/2011GL049945>, 2011.

829 Hatfield, S., Chantry, M., Dueben, P., Lopez, P., Geer, A., and Palmer, T.: Building
830 Tangent-Linear and Adjoint Models for Data Assimilation With Neural Networks, *J. Adv.*
831 *Model. Earth Syst.*, 13, e2021MS002521, <https://doi.org/10.1029/2021MS002521>, 2021.

832 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J.,
833 Nicolas, J., Peubey, C., Radu, R., Schepers, D., Simmons, A., Soci, C., Abdalla, S.,
834 Abellan, X., Balsamo, G., Bechtold, P., Biavati, G., Bidlot, J., Bonavita, M., De Chiara, G.,
835 Dahlgren, P., Dee, D., Diamantakis, M., Dragani, R., Flemming, J., Forbes, R., Fuentes,
836 M., Geer, A., Haimberger, L., Healy, S., Hogan, R. J., Hólm, E., Janisková, M., Keeley, S.,
837 Laloyaux, P., Lopez, P., Lupu, C., Radnoti, G., De Rosnay, P., Rozum, I., Vamborg, F.,
838 Villaume, S., and Thépaut, J.: The ERA5 global reanalysis, *Q. J. R. Meteorol. Soc.*, 146,
839 1999–2049, <https://doi.org/10.1002/qj.3803>, 2020.

840 Hochreiter, S. and Schmidhuber, J.: Long Short-Term Memory, *Neural Comput.*, 9, 1735–
841 1780, <https://doi.org/10.1162/neco.1997.9.8.1735>, 1997.

842 Irrgang, C., Boers, N., Sonnewald, M., Barnes, E. A., Kadow, C., Staneva, J., and
843 Saynisch-Wagner, J.: Towards neural Earth system modelling by integrating artificial
844 intelligence in Earth system science, *Nat. Mach. Intell.*, 3, 667–674,
845 <https://doi.org/10.1038/s42256-021-00374-3>, 2021.

846 James, W. and Stein, C.: Estimation with Quadratic Loss, in: Breakthroughs in Statistics,
847 edited by: Kotz, S. and Johnson, N. L., Springer New York, New York, NY, 443–460,
848 https://doi.org/10.1007/978-1-4612-0919-5_30, 1992.

849 Kadra, A., Lindauer, M., Hutter, F., and Grabocka, J.: Well-tuned Simple Nets Excel on
850 Tabular Datasets, <http://arxiv.org/abs/2106.11189>, 5 November 2021.

851 Keisler, R.: Forecasting Global Weather with Graph Neural Networks,
852 <http://arxiv.org/abs/2202.07575>, 15 February 2022.

853 Kimpson, T., Choulga, M., Chantry, M., Balsamo, G., Boussetta, S., Dueben, P., and
854 Palmer, T.: Deep learning for quality control of surface physiographic fields using
855 satellite Earth observations, *Hydrol. Earth Syst. Sci.*, 27, 4661–4685,
856 <https://doi.org/10.5194/hess-27-4661-2023>, 2023.

857 Kingma, D. P. and Ba, J.: Adam: A Method for Stochastic Optimization, 2017.

858 Kratzert, F., Herrnegger, M., Klotz, D., Hochreiter, S., and Klambauer, G.:
859 NeuralHydrology -- Interpreting LSTMs in Hydrology,
860 <https://doi.org/10.48550/ARXIV.1903.07903>, 2019a.

861 Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., and Nearing, G.: Towards
862 learning universal, regional, and local hydrological behaviors via machine learning
863 applied to large-sample datasets, *Hydrol. Earth Syst. Sci.*, 23, 5089–5110,
864 <https://doi.org/10.5194/hess-23-5089-2019>, 2019b.

865 Lam, R., Sanchez-Gonzalez, A., Willson, M., Wirnsberger, P., Fortunato, M., Alet, F.,
866 Ravuri, S., Ewalds, T., Eaton-Rosen, Z., Hu, W., Merose, A., Hoyer, S., Holland, G.,
867 Vinyals, O., Stott, J., Pritzel, A., Mohamed, S., and Battaglia, P.: Learning skillful
868 medium-range global weather forecasting, *Science*, 382, 1416–1421,
869 <https://doi.org/10.1126/science.adi2336>, 2023.

870 Lang, S., Alexe, M., Chantry, M., Dramsch, J., Pinault, F., Raoult, B., Clare, M. C. A.,
871 Lessig, C., Maier-Gerber, M., Magnusson, L., Bouallègue, Z. B., Nemesio, A. P., Dueben,
872 P. D., Brown, A., Pappenberger, F., and Rabier, F.: AIFS - ECMWF's data-driven
873 forecasting system, <http://arxiv.org/abs/2406.01465>, 3 June 2024.

874 Lees, T., Reece, S., Kratzert, F., Klotz, D., Gauch, M., De Bruijn, J., Kumar Sahu, R., Greve,
875 P., Slater, L., and Dadson, S. J.: Hydrological concept formation inside long short-term
876 memory (LSTM) networks, *Hydrol. Earth Syst. Sci.*, 26, 3079–3101,
877 <https://doi.org/10.5194/hess-26-3079-2022>, 2022.

878 Li, L., Carver, R., Lopez-Gomez, I., Sha, F., and Anderson, J.: SEEDS: Emulation of
879 Weather Forecast Ensembles with Diffusion Models,
880 <https://doi.org/10.48550/ARXIV.2306.14066>, 2023.

881 Machac, D., Reichert, P., and Albert, C.: Emulation of dynamic simulators with
882 application to hydrology, *J. Comput. Phys.*, 313, 352–366,
883 <https://doi.org/10.1016/j.jcp.2016.02.046>, 2016.

- 884 Meyer, D., Grimmond, S., Dueben, P., Hogan, R., and Van Reeuwijk, M.: Machine
885 Learning Emulation of Urban Land Surface Processes, *J. Adv. Model. Earth Syst.*, 14,
886 e2021MS002744, <https://doi.org/10.1029/2021MS002744>, 2022.
- 887 Mironov, D. and Helmert, J.: Parameterization of Lakes in NWP and Climate Models, n.d.
- 888 Muñoz-Sabater, J., Dutra, E., Agustí-Panareda, A., Albergel, C., Arduini, G., Balsamo, G.,
889 Boussetta, S., Choulga, M., Harrigan, S., Hersbach, H., Martens, B., Miralles, D. G.,
890 Piles, M., Rodríguez-Fernández, N. J., Zsoter, E., Buontempo, C., and Thépaut, J.-N.:
891 ERA5-Land: a state-of-the-art global reanalysis dataset for land applications, *Earth*
892 *Syst. Sci. Data*, 13, 4349–4383, <https://doi.org/10.5194/essd-13-4349-2021>, 2021.
- 893 Nath, S., Lejeune, Q., Beusch, L., Seneviratne, S. I., and Schleussner, C.-F.: MESMER-M:
894 an Earth system model emulator for spatially resolved monthly temperature, *Earth Syst.*
895 *Dyn.*, 13, 851–877, <https://doi.org/10.5194/esd-13-851-2022>, 2022.
- 896 Nearing, G., Cohen, D., Dube, V., Gauch, M., Gilon, O., Harrigan, S., Hassidim, A., Klotz,
897 D., Kratzert, F., Metzger, A., Nevo, S., Pappenberger, F., Prudhomme, C., Shalev, G.,
898 Shenzis, S., Tekalign, T. Y., Weitzner, D., and Matias, Y.: Global prediction of extreme
899 floods in ungauged watersheds, *Nature*, 627, 559–563, [https://doi.org/10.1038/s41586-](https://doi.org/10.1038/s41586-900)
900 [024-07145-1](https://doi.org/10.1038/s41586-024-07145-1), 2024.
- 901 Ozaki, Y., Tanigaki, Y., Watanabe, S., Nomura, M., and Onishi, M.: Multiobjective Tree-
902 Structured Parzen Estimator, *J. Artif. Intell. Res.*, 73, 1209–1250,
903 <https://doi.org/10.1613/jair.1.13188>, 2022.
- 904 Pappenberger, F., Ramos, M. H., Cloke, H. L., Wetterhall, F., Alfieri, L., Bogner, K.,
905 Mueller, A., and Salamon, P.: How do I know if my forecasts are better? Using
906 benchmarks in hydrological ensemble prediction, *J. Hydrol.*, 522, 697–713,
907 <https://doi.org/10.1016/j.jhydrol.2015.01.024>, 2015.
- 908 Popov, S., Morozov, S., and Babenko, A.: Neural Oblivious Decision Ensembles for Deep
909 Learning on Tabular Data, <https://doi.org/10.48550/ARXIV.1909.06312>, 2019.
- 910 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., and
911 Prabhat: Deep learning and process understanding for data-driven Earth system
912 science, *Nature*, 566, 195–204, <https://doi.org/10.1038/s41586-019-0912-1>, 2019.
- 913 Sener, O. and Koltun, V.: Multi-Task Learning as Multi-Objective Optimization,
914 <https://doi.org/10.48550/ARXIV.1810.04650>, 2018.
- 915 Shukla, J.: Dynamical Predictability of Monthly Means, *J. Atmospheric Sci.*, 38, 2547–
916 2572, [https://doi.org/10.1175/1520-0469\(1981\)038<2547:DPOMM>2.0.CO;2](https://doi.org/10.1175/1520-0469(1981)038<2547:DPOMM>2.0.CO;2), 1981.
- 917 Shwartz-Ziv, R. and Armon, A.: Tabular Data: Deep Learning is Not All You Need,
918 <https://doi.org/10.48550/ARXIV.2106.03253>, 2021.
- 919 Slater, L. J., Arnal, L., Boucher, M.-A., Chang, A. Y.-Y., Moulds, S., Murphy, C., Nearing,
920 G., Shalev, G., Shen, C., Speight, L., Villarini, G., Wilby, R. L., Wood, A., and Zappa, M.:

- 921 Hybrid forecasting: blending climate predictions with AI models, *Hydrol. Earth Syst.*
922 *Sci.*, 27, 1865–1889, <https://doi.org/10.5194/hess-27-1865-2023>, 2023.
- 923 Thorpe, A., Bauer, P., Magnusson, L., and Richardson, D.: An evaluation of recent
924 performance of ECMWF’s forecasts, <https://doi.org/10.21957/HI1EEKTR>, 2013.
- 925 Van Katwyk, P., Fox-Kemper, B., Seroussi, H., Nowicki, S., and Bergen, K. J.: A Variational
926 LSTM Emulator of Sea Level Contribution From the Antarctic Ice Sheet, *J. Adv. Model.*
927 *Earth Syst.*, 15, e2023MS003899, <https://doi.org/10.1029/2023MS003899>, 2023.
- 928 Viterbo, P.: *Land_surface_processes*, 2002.
- 929 Wesselkamp, M., Moser, N., Kalweit, M., Boedecker, J., and Dormann, C. F.: Process-
930 guidance improves predictive performance of neural networks for carbon turnover in
931 ecosystems, <https://doi.org/10.48550/ARXIV.2209.14229>, 2022.
- 932 Zwart, J. A., Oliver, S. K., Watkins, W. D., Sadler, J. M., Appling, A. P., Corson-Dosch, H.
933 R., Jia, X., Kumar, V., and Read, J. S.: Near-term forecasts of stream temperature using
934 deep learning and data assimilation in support of management decisions, *JAWRA J. Am.*
935 *Water Resour. Assoc.*, 59, 317–337, <https://doi.org/10.1111/1752-1688.13093>, 2023.
- 936