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

Marieke Wesselkamp[1], Matthew Chantry[2], Ewan Pinnington[2], Margarita Choulga[2], Souhail Boussetta[2], Maria Kalweit[3], Joschka Boedecker[3,4], Carsten F. Dormann[1], Florian Pappenberger[2], and Gianpaolo Balsamo[2,5]

1 Department of Biometry, University of Freiburg, Germany
2 European Centre for Medium-Range Weather Forecasts, Reading, United Kingdom
3 Department of Computer Science, University of Freiburg, Germany
4 BrainLinks-BrainTools, University of Freiburg, Germany
5 World Meteorological Organization, Geneva, Switzerland

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

**Abstract**

19

20

21    Most useful weather prediction for the public is near the surface. The processes that are most

22    relevant for near-surface weather prediction are also those that are most interactive and

23    exhibit positive feedback or have key role in energy partitioning. Land surface models

24    (LSMs) consider these processes together with surface heterogeneity and forecast water,

25    carbon and energy fluxes, and coupled with an atmospheric model provide boundary and

26    initial conditions. This numerical parametrization of atmospheric boundaries being

27    computationally expensive, statistical surrogate models are increasingly used to accelerated

28    progress in experimental research. We evaluated the efficiency of three surrogate models in

29    speeding up experimental research by simulating land surface processes, which are integral to

30    forecasting water, carbon, and energy fluxes in coupled atmospheric models. Specifically, we

31    compared the performance of a Long-Short Term Memory (LSTM) encoder-decoder

32    network, extreme gradient boosting, and a feed-forward neural network within a physics-

33    informed multi-objective framework. This framework emulates key states of the ECMWF's

34    Integrated Forecasting System (IFS) land surface scheme, ECLand, across continental and

35    global scales. Our findings indicate that while all models on average demonstrate high

36    accuracy over the forecast period, the LSTM network excels in continental long-range

37    predictions when carefully tuned, the XGB scores consistently high across tasks and the MLP

38    provides an excellent implementation-time-accuracy trade-off. The runtime reduction

39    achieved by the emulators in comparison to the full numerical models are significant, offering

40    a faster, yet reliable alternative for conducting numerical experiments on land surfaces.

41

## 1 Introduction

While forecasting of climate and weather system processes has long been a task for numerical models, the recent development in deep learning has introduced competitive machine-learning (ML) systems for numerical weather prediction (NWP) (Bi et al., 2022; Lam et al., 2023), (Lang et al., 2024). Land surface models (LSMs), even though being an integral part of numerical weather prediction, have not yet caught the attention of the ML-community. LSMs forecast water, carbon and energy fluxes, and in coupling with an atmospheric model, provide the lower boundary and initial conditions [3], [4]. The parametrization of land surface states thus does not only affect predictability of earth and climate systems on sub-seasonal scales (Muñoz-Sabater et al., 2021), but also the short- and medium-range skill of NWP forecasts (De Rosnay et al., 2014). Beyond the online integration with NWPs, offline versions of LSMs provide research tools for experiments on the land surface (Boussetta et al., 2021), the diversity of which are however limited by the required substantial computational resources and often moderate runtime efficiencies (Reichstein et al., 2019).

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

Modelling approaches used for emulation range from low parametrized, auto-regressive linear models to highly non-linear and flexible neural networks (Nath et al., 2022), (Baker et al., 2022), (Chantry et al., 2021), (Meyer et al., 2022). In the global land surface system M-MESMER, a set of simple AR1 regression models is used to initialize the numerical LSM, resulting in a modularized emulator (Nath et al., 2022). Numerical forecasts of gross primary 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 variables from gravity wave drags to

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 & 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 [5] and also constitutes a standalone

87    framework for offline forecasting of land surface processes, the advantage of which for the

88    online framework is the temporal consistency of prognostic state variables (Muñoz-Sabater et

89    al., 2021). The modular construction of ECLand offers potential for element-wise

90    improvement of process representation and thus a stepwise development towards increased

91    computational efficiency. Within the IFS, ECLand also forms the basis of the land surface

92    data assimilation system, updating the land surface state with synoptic data and satellite

93    observations of soil moisture and snow. Emulators of physical systems have been shown to

94    be beneficial in data assimilation routines, allowing for a quick and low maintenance

95    estimation of the tangent linear model (Hatfield et al., 2021). Together with the potential to

96    run large ensembles of land surface states at a much-reduced cost, this would be a potential

97    application of the surrogate models introduced here.

98    Long-short term memory networks (LSTMs) have gained popularity in hydrological

99    forecasting as rainfall-runoff models, for predicting stream flow temperature and also soil

100   moisture [e.g. (Kratzert, Klotz, et al., 2019), (Lees et al., 2022), (Zwart et al., 2023), (Bassi

101   et al., 2024)]. Research on the interpretability of LSTMs has found correlations between the

102   model cell states and spatially or thematically similar hydrological units (Lees et al., 2022),

103   suggesting the specific usefulness of LSTM for representing variables with dynamic storages

104   and reservoirs (Kratzert, Herrnegger, et al., 2019). As emulators, LSTMs have been shown

105   useful for sea surface level projection in a variational manner with Monte Carlo dropout (Van

106   Katwyk et al., 2023). While most of these studies trained their models on observations or

107   reanalysis data, our emulator learns the representation from ECLand simulations directly. To

108    our knowledge, a comparison of models without memory mechanisms to an LSTM-based

109    neural network for 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. These three state variables represent the core of the current

113    configuration of ECLand We specifically focus on the utility of memory mechanisms,

114    highlighting the development of a single LSTM-based encoder-decoder model compared to

115    an extreme gradient boosting approach (XGB) and a multilayer perceptron (MLP), which all

116    perform the same tasks. The LSTM architecture builds on an encoder-decoder network design

117    introduced for flood forecasting (Nearing et al., 2024). To compare forecast skill

118    systematically, the three emulators were compared in long-range forecasting against

119    climatology (Pappenberger et al., 2015). In this work, evaluation is done on ECLand

120    simulation only, i.e. on purely synthetic data, while future work will encompass transfer

121    learning and validation on observations.

122

123    **2 Methods**

124

125    **2.1 The Land Surface Model: ECLand**

126

127    ECLand is a tiled ECMWF Scheme for Surface Exchanges over Land that represents surface

128    heterogeneity and incorporates land surface hydrology (ECLand) (Balsamo et al., 2011)

129    (ECMWF, 2017). ECLand computes surface turbulent fluxes (of heat, moisture and

130    momentum) and skin temperature over different tiles (vegetation, bare soil, snow,

131    interception and water) and then calculates an area-weighted average for the grid-box to

132    couple with the atmosphere (Boussetta et al., 2021). For the overall accuracy of the model,

133    accurate parameterizations are essential (Kimpson et al., 2023) as e.g. the land surface

134    parameterization determines the sensible and latent heat fluxes, and provide the lower

135    boundary conditions for enthalpy and moisture equations in the atmosphere (Viterbo, 2002).

136    We emulate three prognostic state variables of ECLand, that represent core land surface

137    processes: soil water volume and soil temperature at each three depth layers (each at $0-7$

138    cm, $7-21$ cm and $21-72$ cm) and snow cover fraction, aggregated at the surface layer, so

139    below are some more details on these parametrisations.

140

**2.2 Data sources**

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

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

**2.2.1 Surface physiographic fields**

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

**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 & 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 $\sim$ 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 shows the observation-derived static*

194  *physiographic fields, the middle column ERA5 dynamic physiographic and meteorological fields and the rightmost column*

195  *ECLand generated dynamic prognostic state variables.*

| Climate fields | Units | Atmospheric forcing | Units | Prognostic states | Units |
|---|---|---|---|---|---|
| Vegetation cover (*low, high*) | | Total precipitation fraction (*convective + stratiform*) | | Soil water volume (*Layers 1-3*) | m3 m-3 |
| Type of vegetation (*low, high*) | | Downward radiation (*long, short*) | W/m2 | Soil temperature (*Layers 1-3*) | K |
| Minimum stomatal resistance (*low, high*) | | Seasonal LAI (*high, low*) | | Snow cover fraction | |
| Roughness length (*low, high*) | | Wind speed (*v, u*) | m/s | | |
| Urban cover | | Surface pressure | Pa | | |
| Lake cover | | Skin temperature | K | | |
| Lake depth | | | | | |

| Orography (+ *std*, + *filtered*) | m2/s-2 | Specific humidity | kg/kg |
|---|---|---|---|
| Photosynthesis pathways | | Rainfall rate (*total*) | kg/m2s |
| Soil type | | Snowfall rate (*total*) | kg/m2s |
| Glacier mask | | | |
| Permanent wilting point | | | |
| Field capacity | | | |
| Cell area | | | |

196

**2.3 Emulators**

198

199 We compare the utility of a long-short term memory neural network (LSTM), that of extreme

200 gradient boosting regression trees (XGB) and that of a feedforward neural network (that we

201 here refer to as multilayer perceptron, MLP). To motivate this setup and pave the way for

202 discussing effects of (hyper-)parameter choices, a short overview of all approaches is given.

203 All analyses were conducted in Python. XGB was developed in dmlc's XGBoost python

204 package[1]. The MLP and LSTM were developed in the PyTorch lightning framework for deep

205 learning[2]. Neural networks were trained with the Adam algorithm for stochastic optimization

206 (Kingma & Ba, 2017). Model architectures and algorithmic hyperparameters were selected

207 through Bayesian hyperparameter optimization with the Optuna framework (Akiba et al.,

208 2019). The Bayesian optimization minimizes the neural network validation accuracy,

209 specified here as mean absolute error (MAE), over a predefined search space for free

210 hyperparameters with the Tree-structured Parzen Estimator (Ozaki et al., 2022). The resulting

211 hyperparameter and architecture choices which were used for the different approaches are

212 listed in the Supplementary Material.

213

214 **2.3.1 MLP**

215

216 For creation of the MLP emulator we work with a feed-forward neural network architecture

217 of connected hidden layers with ReLU activations and dropout layers, model components

218 which are given in detail in the Supplementary Material or in (Goodfellow et al., 2016). The

---

[1] https://xgboost.readthedocs.io/en/stable/python/index.html
[2] https://lightning.ai/docs/pytorch/stable/

219   MLP was trained with a learning rate scheduler. L2-regularization was added to the training

220   objective via weight decay. Sizes and width of hidden layers as well as hyperparameters were

221   selected together in the hyperparameter optimization procedure. Instead of forecasting

222   absolute prognostic state variables $\boldsymbol{z_t}$, the MLP predicts the 6-hourly increment, $\frac{\widehat{d\boldsymbol{z}}}{dt}$. It is

223   trained on a stepwise rollout prediction of future state variables at a pre-defined lead time at

224   given forcing conditions, see details in the section on optimization.
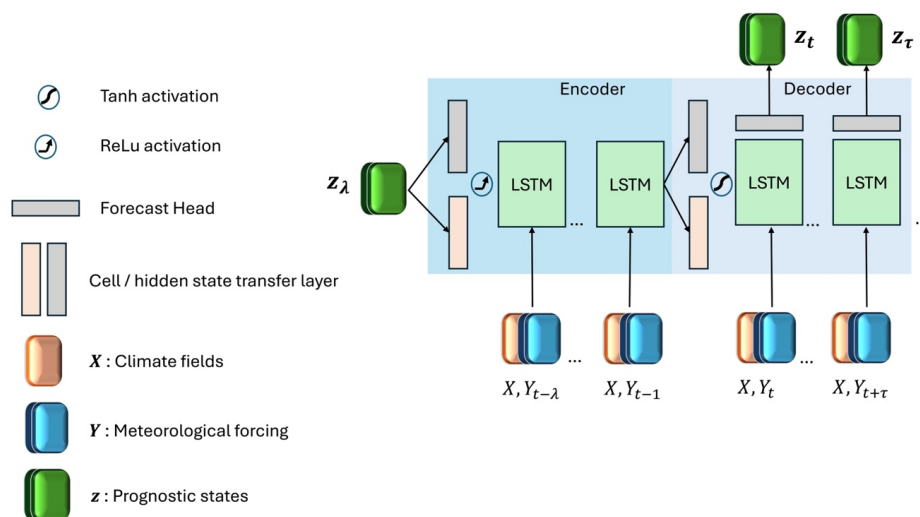
225

226   **2.3.2 LSTM**

227

228   LSTMs are recurrent networks that consider long-term dependencies in time series through

229   gated units with input and forget mechanisms (Hochreiter & Schmidhuber, 1997). In

230   explicitly providing time-varying forcing and state variables, LSTM cell states serve as long-

231   term memory while LSTM hidden states are the cells' output and pass on stepwise short-term

232   representations stepwise. In short notation (Lees et al., 2022), a one-step ahead forward pass

233   followed by a linear transformation can be formulated as

234   $$\boldsymbol{h_t}, \boldsymbol{c_t} = f(\boldsymbol{x_t}, \boldsymbol{h_{t-1}}, \boldsymbol{c_{t-1}}, \boldsymbol{\theta})$$

235   $$\hat{\boldsymbol{z}}_t = \boldsymbol{A}\boldsymbol{h_t} + b$$

236   where $\boldsymbol{h_{t-1}}$ denotes the hidden state, i.e. output estimates from the previous time step, $\boldsymbol{c_{t-1}}$

237   the cell state from the previous time step, and $\boldsymbol{\theta}$ the time-invariant model weights. We stacked

238   multiple LSTM cells to an encoder-decoder model with transfer layers for hidden and cell

239   state initialization and for transfer to the context vector (see figure 1) (Nearing et al., 2024). A

240   lookback $l$ of the previous static and dynamic feature states are passed sequentially to the first

241   LSTM cells in the encoder layer, while the $l$ prognostic state variables $\boldsymbol{z}$ initialize the hidden

242   state $\boldsymbol{h}_0$ after a linear embedding. The output of the first LSTM layer cells become the input

243   to the deeper LSTM layer cells and the last hidden state estimates are the final output from

244   the encoder. Followed by a non-linear transformation with hyperbolic tangent activation, the

245   hidden cell states are transformed into a weighted context vector $\boldsymbol{s}$. Together with the encoder

246   the cell state $(\boldsymbol{c_t}, \boldsymbol{s})$ initializes the hidden and cell states of the decoder. The decoder LSTM

247   cells take as input again static and dynamic features sequentially at lead times $t = 1, \dots, \tau$, but

248   not the prognostic states variables. These are estimated from the sequential hidden states of

249   the last LSTM layer cells, transformed to target size with a linear forecast head before

250   prediction. LSTM predicts absolute state variables $\boldsymbol{z_t}$ while being optimized on $\boldsymbol{z_t}$ and $d\hat{\boldsymbol{z}}_t$

251   simultaneously, see section on optimization.

252
253 *Figure 1: LSTM architecture. Blue shaded area indicates the encoder part, where the model is driven by a lookback λ*
254 *of meteorological forcing and state variables. The light-blue shaded area indicates the decoder part that is initialized*
255 *from the encoding to unroll LSTM forecasts from the initial time step t up to a flexibly long lead time of τ.*

256 **2.3.3 XGB**

257

258 Extreme gradient boosting (XGB) is a regression tree ensemble method that uses an

259 approximate algorithm for best split finding. It computes first and second order gradient

260 statistics in the cost function, performing a similar to gradient descent optimization (T. Chen

261 & Guestrin, 2016), where each new learner is trained on the residuals of the previous ones.

262 Regularization and column sampling aim for preventing overfitting internally. XGB is known

263 to provide a powerful benchmark for time series forecasting and tabular data [(T. Chen &

264 Guestrin, 2016; Shwartz-Ziv & Armon, 2021), (X. Chen et al., 2020)]. Like the MLP, it is

265 trained to predict the increment $\widehat{dz}_{t,i}$ of prognostic state variables, but only for a one-step

266 ahead prediction.

267

268 **2.4 Experimental setup**

269

270 We distinguish the experimental analysis into three parts that vary in the usage of the training

271 database: (1) model development, (2) model testing, and (3) global model transfer.

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

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

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

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

276    selection and model performances, see Supplementary Material S3.

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

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

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

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

281    2022 were used, see details in section on forecast horizons. With this same data splitting

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

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

284    a systematic comparison of the three models, because high resolution training with XGB was

285    prohibited by the required working memory. On the other hand, this extrapolation scenario

286    created an unseen problem for the models that were selected on a continental and high-

287    resolution scale which is reflected in the resulting scores.

288

289    **2.5 Optimization**

290

291    **2.5.1 Loss functions**

292

293    The basis of the loss function $\mathcal{L}$ for the neural network optimization was PyTorch's

294    SmoothL1Loss[3], a robust loss function that combines L1-norm and L2-norm and is less

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

296    parameter $\beta$, smooth L1 transitions from L2-norm to L1-norm above the threshold.

297    SmoothL1Loss $\mathcal{L}$ is defined as

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

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

300    here with $\beta = 1$. All models were trained to minimize the incremental loss $\mathcal{L}_s$ that is the

301    differences between the estimates of the seven prognostic states increments $\widehat{d\mathbf{z}_t}$ and the full

302    model's prognostic states increments $d\mathbf{z}_t$ simultaneously as the sum of losses over all states.

303    We opted for a loss function equally weighted by variables to share inductive biases among

304    the non-independent prognostic states (Sener & Koltun, 2018). When aggregating over all

305    training lead times $t = 1, \dots, \tau$, $\mathcal{L}_s$ and grid cells $i = 1, \dots, p$ is

---

[3] https://pytorch.org/docs/stable/generated/torch.nn.SmoothL1Loss.html

306
$$\mathcal{L}_s(\widehat{d\mathbf{z}}, d\mathbf{z}) = \sum_t^{\tau} \sum_i^{p} \mathcal{L}_t(\widehat{d\mathbf{z}}_{t,i}, d\mathbf{z}_{t,i}),$$

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

308

309
$$\mathcal{L}_r(\widehat{d\mathbf{z}}, \mathbf{z}) = \frac{1}{\tau} \sum_t^{\tau} \sum_i^{p} \mathcal{L}_t(z_{t-1,i} + \widehat{d\mathbf{z}}_{t,i}, z_{t,i})$$

310

311  Prognostic state increments are essentially the first differences from one to the next timestep

312  that are normalized again by the global standard deviation of the model's states increments,

313  $s_{dz}$ before computation of the loss (Keisler, 2022). Due to the forecast models' structural

314  differences, loss functions were individually adapted:

315  **MLP** The combined loss function for the MLP is the sum of the incremental loss $\mathcal{L}_s$ and the

316  rollout loss $\mathcal{L}_r$. For the rollout loss $\mathcal{L}_r$, $\mathcal{L}$ was aggregated over grid cells $p$ and accumulated

317  after an auto-regressive rollout over lead times $\tau$, before being averaged out by division by $\tau$

318  (Keisler, 2022).

319  **LSTM** The combined loss function for the LSTM is the sum of the incremental loss

320  $\mathcal{L}_s$, where the $d\hat{\mathbf{z}}_t$ were derived from $\hat{\mathbf{z}}_t$ after the forward pass, and the loss $\mathcal{L}$ computed on

321  decoder estimates of prognostic states variables, a functionality that leverages the potential of

322  our LSTM structure.

323  **XGB** Trained only from one to the next time step, i.e. at a lead time of $\tau = 1$, the incremental

324  loss $\mathcal{L}_s = \mathcal{L}_r$. Without a SmoothL1Loss implementation provided in dmlc's XGBoost, we

325  trained XGB with both the Huber-Loss and the default L2-loss. The latter initially providing

326  better results, we chose the default L2-norm as loss function for XGB with the regularization

327  parameter $\lambda = 1$.

328

329  **2.5.1 Normalization**

330  As prognostic target variables are all lower bounded by zero, we tested both z-scoring and

331  max-scoring. The latter yielded no significant improvement, thus we show our results with z-

332  scored target variables. For neural network training but not for fitting XGB, static, dynamic

333  and prognostic state variables were all normalized with z-scoring towards the continental or

334  global mean $\bar{z}$ and unit standard deviation $s_z$ as

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

336    Prognostic target state increments were normalized again by the global standard deviation of

337    increments computing the loss (see section 2.5.1) to smooth magnitudes of increments

338    (Keisler, 2022). State variables were backtransformed to original scale before evaluation.

339

340    **2.5.3 Spatial and temporal sampling**

341    Sequences were sampled randomly from the training data set, while validation happened

342    sequentially. MLP and XGB were trained on all grid cells simultaneously in both the

343    continental and global setting, while LSTM was trained on the full continental data set but

344    was limited by GPU memory in the global task. We overcame this limitation by randomly

345    subsetting grid cells in the training data into largest possible, equally sized subsets which

346    were then loaded along with the temporal sequences during the batch sampling.

347

348    **2.6 Evaluation**

349

350    Three scores are used for model validation during the model development phase and in

351    validating architecture and hyperparameter selection, being the root mean squared error

352    ($RMSE$), the mean absolute error ($MAE$) and the anomaly correlation coefficient ($ACC$).

353    First, scores were assessed objectively in quantifying forecast accuracy of the emulators

354    against ECLand simulations directly with RMSE and MAE. Doing so, scores were

355    aggregated over lead times, grid cells or both. The total RMSE was computed as

356
$$\text{RMSE} = \sqrt{\frac{\sum_{\tau,p}(z - \hat{z})^2}{n}},$$

357    As the mean absolute error in prognostic state variable prediction over the total of $n$ grid cells

358    $p$ times lead times $\tau$. Equivalently, MAE was computed as

359
$$\text{MAE} = \frac{\sum_{t,p}|z - \hat{z}|}{n},$$

360    Beyond accuracy, the forecast skill of emulators was assessed using a benchmark model: the

361    ACC (see below) as index of the long-term naïve climatology $c$ of ECLand, forced by ERA5

362    (see section 2.2). More specifically, this is the 6-hourly mean of prognostic state variables

363    over the last 10 years preceding the test year, i.e. the years 2010 to 2020. While climatology

364    is a hard-to-beat benchmark specifically in long-term forecasting, the persistence is a

365    benchmark for short-term forecasting (Pappenberger et al., 2015). For verification against

366    climatology, we compute the anomaly correlation coefficient (ACC) over lead times as

367
$$\text{ACC}(t) = \frac{\overline{(\hat{z} - c)(z - c)}}{\sqrt{\overline{(\hat{z} - c)^2}\ \overline{(z - c)^2}}}$$

368    at each t = 1, ..., $\tau$ where the overbar denotes averaging over grid cells $p = i, ..., n$. The

369    nominator now indicates the mean squared skill error towards climatology and the

370    denominator its variability. ACC is bounded between 1 and -1, and an ACC of 1 indicates

371    perfect representation of forecast error variability, an ACC of 0.5 indicates a similar forecast

372    error to that of the climatology, an ACC of 0 indicates that forecast error variability

373    dominates and the forecast has no value and an ACC approaching -1 indicates that the

374    forecast has been very unreliable (ECMWF, n.d.). ACC is undefined when the denominator

375    is zero. This is the case either when mean squared emulator or ECLand anomaly, or both are

376    zero because forecast and climatology perfectly align, or because they cancel out at

377    summation to the mean.

378

379    **2.6.1 Forecast horizons**

380

381    Forecast horizons of the emulators are defined by the decomposition of the RMSE

382    (Bengtsson et al., 2008) into the emulator's variability around climatology (i.e. anomaly),

383    ECLand's variability around climatology and the covariance of both. The horizon is the point

384    in time at which the forecast error reaches saturation level, that is when the covariance of

385    emulator and ECLand anomalies approaches zero, as does the ACC.

386    We analysed predictive ability and predictability by computing the ACC for all lead times

387    from 6 hours to approx. one year, i.e. lead times $t = 1, ..., \tau$, $\tau$ being 1350. As this confounds

388    the seasonality with the lead time, we compute these for every starting point of the prediction,

389    requiring two test years (2021 and 2022).

390    Forecast horizons based on the emulators' skill in standardized anomaly towards persistence

391    were equivalently computed but with persistence as a benchmark for shorter time scales, this

392    was only done for three months, from January to March 2021.

393    The analysis was conducted on two exemplary regions in northern and southern Europe that

394    represent very different conditions orography and in prognostic land surface states,

395    specifically in snow cover. For details on the regions and on the horizons computed with

396    standardized anomaly skill, see Appendices A1 and A4 respectively.

397

398    **3 Results**

### 3.1 Aggregated performances

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

At variable level, results differentiate into model specific strengths. In soil water volume, XGB outperforms the neural network emulators by up to 60% in the first and second layer MAEs towards the LSTM and up to over 40% for towards the MLP (see table 3). While the representation of anomalies by specifically the LSTM decreases towards lower soil layers with an ACC of only 0.6214 at the third soil layer, it remains consistently higher for XGB with an ACC still > 0.789 at soil layer three.

In soil temperature approximation, LSTM achieves best accuracies at higher soil levels with up to 7% improvement in MAE towards XGB and ACCs > 0.92, but XGB outperforms LSTM at the third soil level with a close to 50% improvement (see table 4). The MLP doesn't stand out by high scores on the continental scale. However, in terms of accuracy we found an inverse ranking in the model development procedure during which LSTM outscored XGB in soil water volume but struggled with soil temperature approximations, for the interested reader we refer to the supplementary information.

In snow cover approximation, the LSTM emulator enhances accuracies by over ~50% in MAE towards both the XGB and the MLP emulator and scores highest in anomaly representation with an ACC of ~0.87 compared to an ACC of ~0.66 for the MLP and only ~0.74 for the XGB (see table 5).

**Globe.** Score ranking on the global scale varies strongly from the continental scale (see table 2). In total average accuracies, the MLP outperforms XGB by over 30% and LSTM by up ~25% in RMSE and improves MAE more than 15% towards both. In anomaly correlation however it scores last, whereas XGB achieves the highest total average of over 0.75.

Consistent with scores on the continental scale is XGBs high performance in soil temperature (see table 3). It significantly outperforms the LSTM by ~60% in RMSE and nearly up to 75% in MAE in all layers and the MLP by up to 50% in MAE at the top layer. Anomaly

433    persistence for all models degrade visibly towards the lower soil layers, while that of the

434    LSTM most relative to MLP and XGB. Similar to the continental scale, XGB also

435    outperforms the other candidates in soil temperature forecasts in all but the medium layer,

436    where the MLP gets higher scores in MAE and RMSE but not in ACC (see table 4). LSTM

437    doesn't stand out with any scores on the global scale.

438

439    **3.2 Spatial and temporal performances**
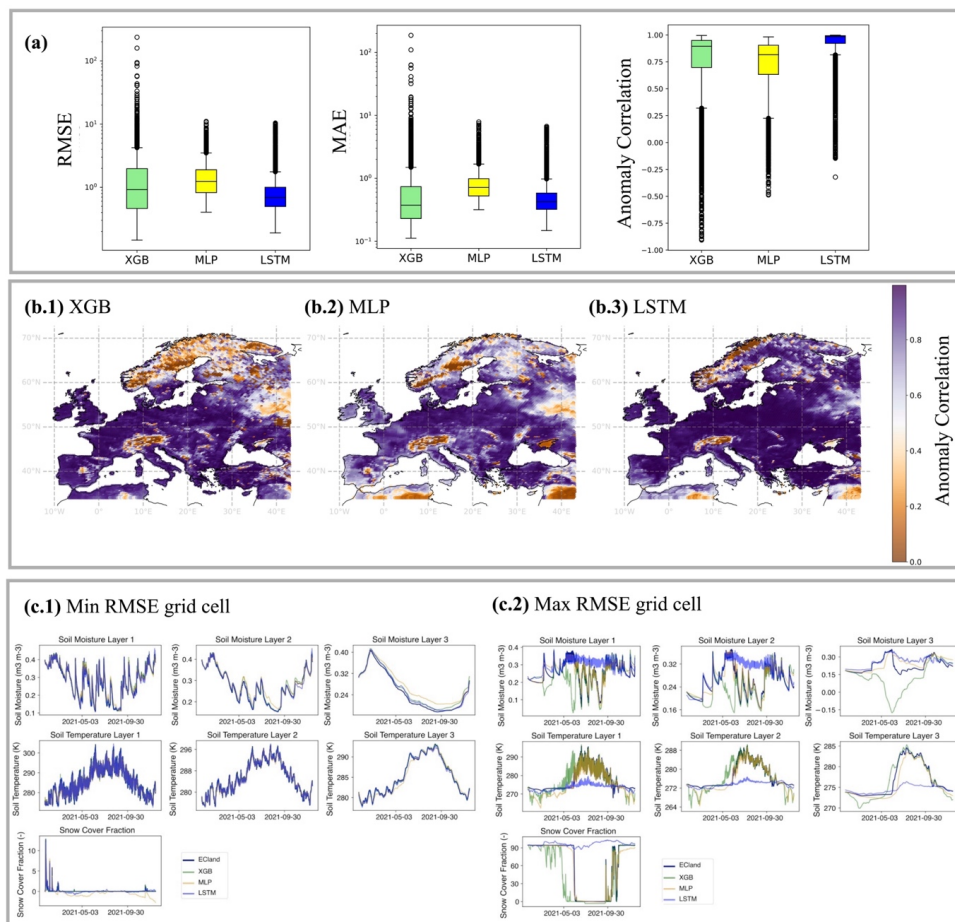
440

441    **Europe.** When summarizing temporally aggregated scores as boxplots to a total distribution

442    over space (see figure 2, A), the long tails of XGB scores become visible, whereas the MLP

443    indicates most robustness. This is reflected in the geographic distribution of scores at the

444    example of ACC (see figure 2, bottom), where the area of low anomaly correlation is largest

445    for XGB, ranging over nearly all northern Scandinavia, while MLP and LSTM have smaller

446    and more segregated areas of clearly low anomaly correlation. The LSTM shows a

447    homogenously high ACCs over most of central Europe but the Alps, while also seems to be

448    challenged in areas of relative to the central Europe extreme weather conditions at the

449    Norwegian and Spanish coasts.

450    **Globe.** Similar to the results from the continental analysis, we find again long upper tails of

451    outliers for XGB in total spatial distribution of accuracies, both in RMSE and MAE and only

452    few outliers for MLP and LSTM. The anomaly correlation distribution changed towards

453    longer lower tails for MLP and LSTM and a shorter lower tail for XGB. We should, however,

454    take the results of total average ACC with care as it remains largely undefined in regions

455    without much noise in snow cover or soil water volume and globally represents mainly

456     patterns of soil temperature.



457

*Figure 2: **a**: Total aggregated distributions of (log) scores averaged over lead times, i.e. displaying the variation among*
458     *grid cells. **b**: The distribution of the anomaly correlation in space on the European subset (b.1: XGB, b.2: MLP, b.3:*
459     *LSTM). **c**: Model forecasts over test year 2021 for grid cell with minimum and maximum RMSE values (LSTM).*
460

461

462     *Table 2: Emulator total average scores, aggregated over variables, time and space from the European and Global*
463     *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 |

Table 3: Emulator average scores on soil water volume forecasts for the European subset, aggregated over space and time from the European and Global 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 |

Table 4: Emulators' mean scores on soil temperature forecasts for the 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 |

Table 5: Emulators' mean scores on snow cover forecasts for the European subset, aggregated over space and time.

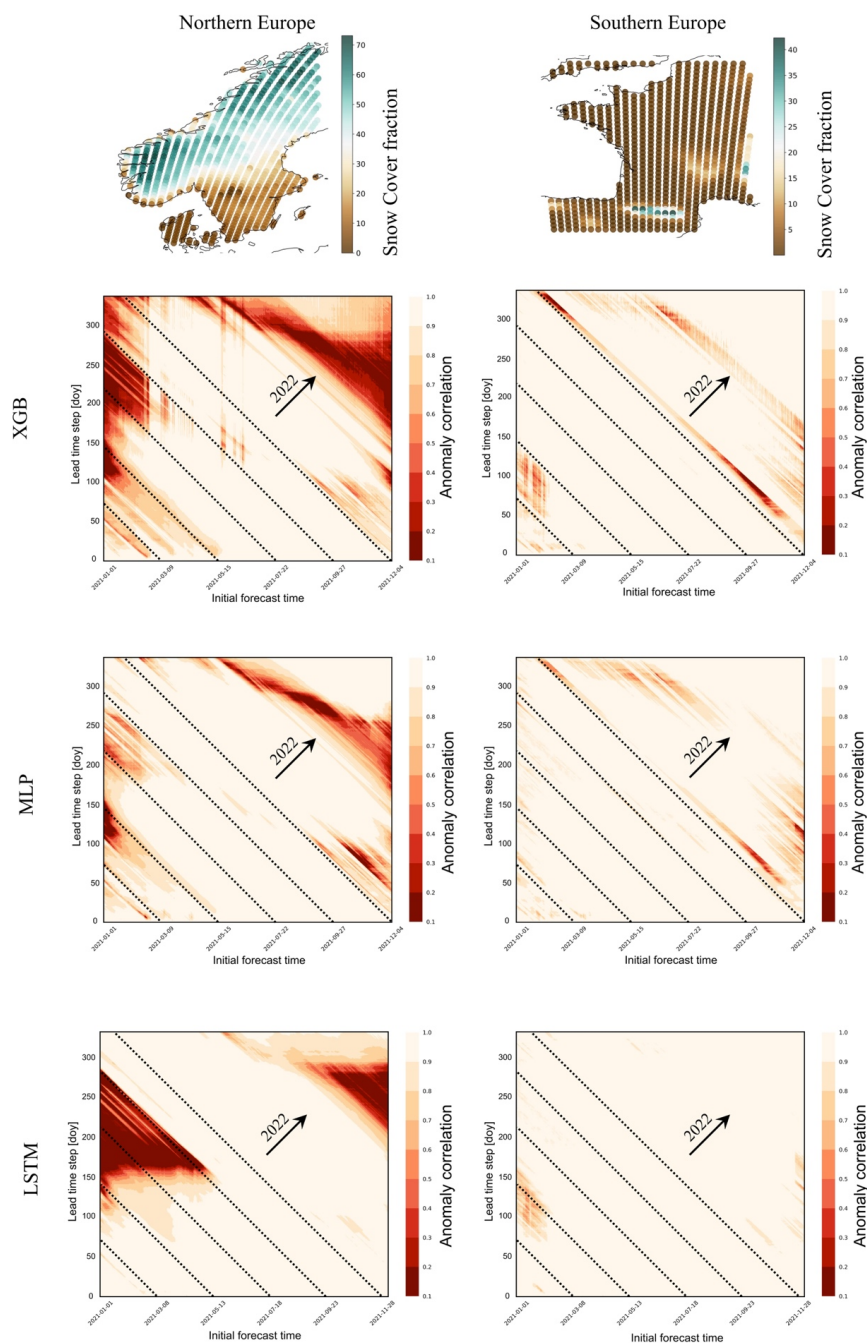| Variable | Layer | Model | RMSE | | MAE | | ACC | |
|---|---|---|---|---|---|---|---|---|
| | | | Europe | Globe | Europe | Globe | Europe | Globe |
| Snow cover | top | XGB | 8.219 | 9.906 | 3.099 | 5.196 | 0.746 | **0.707** |
| | | MLP | 6.449 | **5.995** | 2.986 | **3.671** | 0.66 | 0.618 |
| | | LSTM | **3.526** | 6.127 | **1.47** | 4.357 | **0.877** | 0.698 |

**3.3 Forecast horizons**

473    Forecast horizons were computed for two European regions, of which the northern one
474    represents the area of lowest emulators' skill (see figure 2, B.1-3) and the southern one an
475    area of stronger emulators' skill. Being strongly correlated with soil water volume, these two
476    regions differ specifically in their average snow cover fraction (see figure 3). The displayed
477    horizons were computed over all prognostic state variables simultaneously, while their
478    interpretation is related to horizons computed for prognostic state variables separately, for the
479    figures of which we refer to the Supplementary Material.
480    In the North, predictive skill depended on an interaction of how far ahead a prediction was
481    made (the lead time) and the day of year to which the prediction was made. In the best case,
482    the LSTM, summer predictions were poor (light patches in figure 3 heat maps), but only
483    when initialised in winter. Or, in other words, one can make good predictions starting in
484    winter, but not to summer. Vertical structures indicate a systematic model error that appears at
485    specific initialisation times and that is independent of prediction date, for example in XGB
486    forecasts that are initialized in May (see figure 3, northern region). Diagonal light structures
487    in the heat maps indicate a temporally consistent error and can be interpreted as physical
488    limits of system predictability, where the different initial forecast time doesn't affect model
489    scores.
490    All models show stronger limits in predictability and predictive ability in the northern
491    European region (see figure 3, left column). MLP and XGB struggled with representing
492    seasonal variation towards climatology at long lead times, while LSTM is strongly limited by
493    a systematic error in certain regions. Initializing the forecast the 1 January 2021, MLP drops
494    below an ACC of 80% repeatedly from initialization on and then to an ACC below 10% at the
495    beginning of May. LSTMs performance is more robust in the beginning of the year but
496    depletes strongly later to less than 10% ACC in mid May. On the one hand, this represents
497    two different characteristics of model errors: MLP forecasts for snow cover fraction are less
498    than zero for some grid cells while LSTM forecasts for snow cover fraction remain falsely at
499    very high levels for some grid cells, not predicting the snowmelt in May (see Supplementary
500    Material, S4.1). On the other hand, this represents a characteristic error due to change in
501    seasonality: the snowmelt in this region in May happens abruptly and all emulators
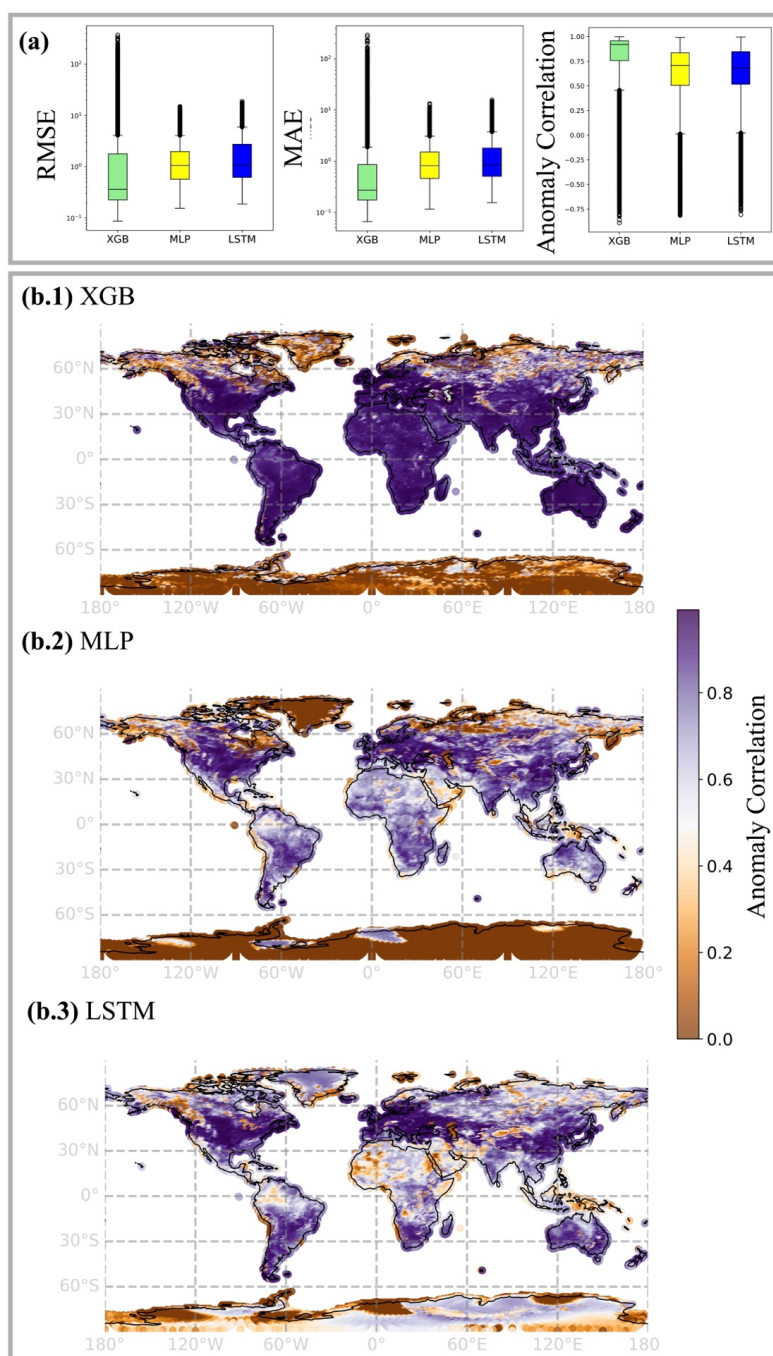502    repeatedly over- or underpredict the exact date.
503

504

*Figure 3: Emulator forecast skill horizons in two European subregions, aggregated over prognostic state variables. Scores are computed with the anomaly correlation coefficient (ACC) at 6-hourly lead times (y-axis) over approx. one year, displayed as a function of the initial forecast time (x-axis). As horizon we define the time at which the forecast has no value at all, i.e. when ACC is 0 (or below 10%). The diagonal dashed lines indicate the day of the test year 2021 as labelled on the x-axis, the arrows indicate where forecasts reach the second test year 2022.*

510



511

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

## 4 Discussion

In the comparative analysis of emulation approaches for land surface forecasting, three primary models—LSTM (Long Short-Term Memory networks), MLP (Multi-Layer Perceptrons), and XGB (Extreme Gradient Boosting)—have been evaluated to understand their effectiveness across different operational scenarios. While all emulators achieved high predictive scores, models differ in their demand of computational resources (Cui et al., 2021) and each one offers unique advantages and faces distinct challenges, impacting their suitability for various forecasting tasks. With this work we want to present the first steps towards enabling quick offline experimentation on the land surface with ECMWF's land surface scheme ECLand and decreasing computational demands, i.e. in the coupled data assimilation.

### 4.1 Approximation of prognostic land surface states

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

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

548 Inversely to expectations and preceding experiments, on the European data set relative to the

549 two other models the LSTM scored better in the upper layer soil temperatures than in

550 forecasting soil water volume and decreased in scores towards lower layers with slower

551 processes. For training on observations, the decreasing LSTM predictive accuracy for soil

552 moisture with lead time is discussed (Datta & Faroughi, 2023), but reasons arising from the

553 engineering side remain unclear. In an exemplary case of a single-objective, deterministic

554 streamflow forecast, a decrease in recurrent neural network performance has been related

555 with an increasing coefficient of variation (Guo et al., 2021). In our European subregions, the

556 signal-to-noise ratio of the prognostic state variables (computed as the averaged ratio of mean

557 and standard deviation) is up to ten times higher in soil temperature than in soil water volume

558 states (see Supplementary Material, S2.1). While a small signal of the latter may induce

559 instability in scores, it does not explain the decreasing performance towards deeper soil layers

560 with slow processes, where we expected an advantage of the long-term memory.

561 Stein's paradox tells us that joint optimization may lead to better results if the target is multi-

562 objective, but not if we are interested in single targets (James & Stein, 1992)(Sener & Koltun,

563 2018). While from a process perspective multi-objective scores are less meaningful than

564 single ones, this is what we opted for due to efficiency. The unweighted linear loss

565 combination might be suboptimal in finding effective parameters across all prognostic state

566 variables (Z. Chen et al., 2017)(Sener & Koltun, 2018), yet being strongly correlated, we

567 deemed their manual weighting inappropriate. An alternative to this provides adaptive loss

568 weighting with gradient normalisation (Z. Chen et al., 2017).

569

**570** **4.2 Evaluation in time and space**

571 We used aggerated MAE and RMSE accuracies as a first assessment tool to conduct model

572 comparison, but score aggregation hides model specific spatio-temporal residual patterns.

573 Further, both scores are variance dependent, favouring low variability in model forecasts

574 even though this may not be representative of the system dynamic (Thorpe et al., 2013).

575 Assessing the forecast skill over time as the relative proximity to a subjectively chosen

576 benchmark helps disentangling areas of strengths and weaknesses in forecasting with the

577 emulators (Pappenberger et al., 2015). The naïve 6-hourly climatology as benchmark

578 highlights periods where emulators long-range forecasts on the test year are externally limited

579 by seasonality, i.e. system predictability, and where they are internally limited by model error,

580 i.e. the model's predictive ability. Applying this strategy in two exemplary European

581 subregions showed that all emulators struggle most in forecasting the period from late

582    summer to autumn, unless they are initialized in summer (see figure 3). Because forecast

583    quality is most strongly limited by snow cover (see Supplementary Material, A4.1), we

584    interpret this as the unpredictable start of snow fall in autumn. External predictability

585    limitations seem to affect the LSTM overall less than the two other models, and specifically

586    XGB drifts at long lead times.

587    From a geographical perspective inferred from the continental scale, emulators struggle in

588    forecasting prognostic state variables in regions with complicated orography and strong

589    environmental gradients. XGB scores vary seemingly random in space, while neural

590    networks scores exhibit spatial autocorrelation. A meaningful inference about this, however,

591    can only be conducted in assessing model sensitivities to physiographic and meteorological

592    fields through gradients and partial dependencies. While the goal of this work is to introduce

593    our approach to emulator development, we envision this for follow-up analyses.

594

595    **4.3 Emulation with memory mechanisms**

596

597    Without much tuning, XGB challenges both LSTM and MLP for nearly all variables (see

598    tables 2-4). In training on observations for daily short-term and real-time rainfall-runoff

599    prediction, XGB and LightXGB were shown before to equally performed as, or outperformed

600    LSTMs (X. Chen et al., 2020)(Cui et al., 2021). Nevertheless, models with memory

601    mechanism such as the encoder-decoder LSTM remain a promising approach for land surface

602    forecasting regarding their differentiability (Hatfield et al., 2021), their flexible extension of

603    lead times, for exploring the effect of long-term dependencies or for inference from the

604    context vector that may help identifying the process relevant climate fields (Lees et al.,

605    2022).

606    In our LSTM architecture, we assume that our model is well defined in that the context vector

607    perfectly informs the hidden decoder states. If that assumption is violated, potential strategies

608    are to create a skip-connection between context vector and forecast head, or to consider input

609    of time-lagged variables or self-attention mechanisms (X. Chen et al., 2020). With attention,

610    the context vector becomes a weighted sum of alignments that relates neighbouring positions

611    of a sequence, a feature that could be leveraged for forecasting quick processes such as snow

612    cover or top-level soil water volume.

613    Comparing average predictive accuracies across different training lead times indicates that

614    training at longer lead times may enhance short-term accuracy of the LSTM at the cost of

615    training runtime (see Supplementary Material, S2). A superficial exploration of encoder

616     length indicates no visible improvement on target accuracies if not a positive tendency

617     towards shorter sequences. This needs an extended analysis for understanding, yet without a

618     significant improvement by increased sequence length, GRU cells might provide a simplified

619     and less parameterized alternative to LSTM cells. They were found to perform equally well

620     on streamflow forecast performance before, while reaching higher operational speed (Guo et

621     al., 2021).

622

623     **4.4 Emulators in application**

624

625     LSTM networks with a decoder structure are valued for their flexible and fast lead time

626     evaluation, which is crucial in applications where forecast intervals are not consistent. The

627     structure of LSTM is well-suited for handling sequential data, allowing it to perform

628     effectively over different temporal scales (Hochreiter & Schmidhuber, 1997). They provide

629     access to gradients, which facilitates inference, optimization and usage for coupled data

630     assimilation (Hatfield et al., 2021). Nevertheless, the complexity of LSTMs introduces

631     disadvantages: Despite their high evaluation speed and accuracy under certain conditions,

632     they require significant computational resources and long training times. They are also highly

633     sensitive to hyperparameters, making them challenging to tune and slow to train, especially

634     with large datasets.

635     MLP models stand out for their implementation, training and evaluation speed with yet

636     rewarding accuracy, making them a favourable choice for scenarios that require rapid model

637     deployment. They are tractable and easy to handle, with a straightforward setup that is less

638     demanding computationally than more complex models. MLPs also allow for access to

639     gradients, aiding in incremental improvements during training and quick inference (Hatfield

640     et al., 2021). Despite these advantages, MLPs face challenges with memory scaling during

641     training at fixed lead times, which can hinder their applicability in large-scale or high-

642     resolution forecasting tasks.

643     XGB models are highly regarded for their robust performance with minimal tuning,

644     achieving high accuracy not only in sample applications, but also in transfer to unseen

645     problems (Shwartz-Ziv & Armon, 2021) (Grinsztajn et al., 2022). Their simplicity makes

646     them easy to handle, even for users with limited technical expertise in machine learning.

647     However, the slow evaluation speed of XGB becomes apparent as dataset complexity and

648     size increase. Although generally more interpretable than deep machine learning tools, XGB

649     is not differentiable, limiting its application in coupled data assimilation (Hatfield et al.,

650     2021) even though research on differentiable trees is ongoing (Popov et al., 2019).

651

652     **5 Conclusion**

653

654     In conclusion, the choice between LSTM, MLP, and XGB models for land surface forecasting

655     depends largely on the specific requirements of the application, including the need for speed,

656     accuracy, and ease of use. Each model's computational demands, flexibility, and operational

657     overhead must be carefully considered to optimize performance and applicability in diverse

658     forecasting environments. When it comes to accuracy, combined model ensembles of XGB

659     and neural networks have been shown to yield the best results (Shwartz-Ziv & Armon, 2021),

660     but accuracy alone will not determine a single best approach (Bouthillier et al., 2021). Our

661     comparative assessment underscores the importance of selecting the appropriate emulation

662     approach based on a clear understanding of each model's strengths and limitations in relation

663     to the forecasting tasks at hand. By developing the emulators for ECMWF's numerical land

664     surface scheme ECLand, we path the way towards a physics-informed ML-based land surface

665     model that on the long run can be parametrized with observations and provide a pretrained

666     model suite to improve land surface forecasts.

667

668     **Code and data availability**

669     Code for this analysis can be found here: https://github.com/MWesselkamp/land-surface-

670     emulation. Data is available on request.

671     **Author contribution**

672     MW, MCha, EP, FP and GB conceived the study. MW and EP conducted the analysis. MW,

673     MCha, MK, EP discussed and took technical decisions. SB advised on process decisions.

674     MW, MCho and FP wrote the manuscript. MW, MCha, EP, MCho, SB, MK, CFD, FP

675     reviewed the analysis and/or manuscript.

676     **Competing interest**

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

678     **Acknowledgements**

679     This work profited from discussion with Linus Magnusson, Sina R. K. Farhadi and Karan

680     Ruparell and many more. MW thankfully acknowledges ECMWF for providing two research

681     visit stipendiates over the course of the collaboration. ChatGPT version 4.0 was used for

682     coding support.

683

**References**

685

686     Akiba, T., Sano, S., Yanase, T., Ohta, T., & Koyama, M. (2019). Optuna: A Next-generation

687          Hyperparameter Optimization Framework. *Proceedings of the 25th ACM SIGKDD*

688          *International Conference on Knowledge Discovery & Data Mining*, 2623–2631.

689          https://doi.org/10.1145/3292500.3330701

690     Baker, E., Harper, A. B., Williamson, D., & Challenor, P. (2022). Emulation of high-

691          resolution land surface models using sparse Gaussian processes with

692          application to JULES. *Geoscientific Model Development*, *15*(5), 1913–1929.

693          https://doi.org/10.5194/gmd-15-1913-2022

694     Balsamo, G., Boussetta, S., Dutra, E., Beljaars, A., & Viterbo, P. (2011). *Evolution of land-*

695          *surface processes in the IFS*. 127.

696     Bassi, A., Höge, M., Mira, A., Fenicia, F., & Albert, C. (2024). *Learning Landscape*

697          *Features from Streamflow with Autoencoders*. https://doi.org/10.5194/hess-

698          2024-47

699     Bengtsson, L. K., Magnusson, L., & Källén, E. (2008). Independent Estimations of the

700          Asymptotic Variability in an Ensemble Forecast System. *Monthly Weather*

701          *Review*, *136*(11), 4105–4112. https://doi.org/10.1175/2008MWR2526.1

702     Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X., & Tian, Q. (2022). *Pangu-Weather: A 3D High-*

703          *Resolution Model for Fast and Accurate Global Weather Forecast*.

704          https://doi.org/10.48550/ARXIV.2211.02556

705    Boussetta, S., Balsamo, G., Arduini, G., Dutra, E., McNorton, J., Choulga, M., Agustí-

706         Panareda, A., Beljaars, A., Wedi, N., Munõz-Sabater, J., De Rosnay, P., Sandu, I.,

707         Hadade, I., Carver, G., Mazzetti, C., Prudhomme, C., Yamazaki, D., & Zsoter, E.

708         (2021). ECLand: The ECMWF land surface modelling system. *Atmosphere*, *12*(6),

709         723. https://doi.org/10.3390/atmos12060723

710    Bouthillier, X., Delaunay, P., Bronzi, M., Trofimov, A., Nichyporuk, B., Szeto, J., Sepah, N.,

711         Raff, E., Madan, K., Voleti, V., Kahou, S. E., Michalski, V., Serdyuk, D., Arbel, T.,

712         Pal, C., Varoquaux, G., & Vincent, P. (2021). *Accounting for Variance in Machine*

713         *Learning Benchmarks* (arXiv:2103.03098). arXiv. http://arxiv.org/abs/2103.03098

714    Chantry, M., Hatfield, S., Duben, P., Polichtchouk, I., & Palmer, T. (2021). *Machine*

715         *learning emulation of gravity wave drag in numerical weather forecasting*.

716         https://doi.org/10.48550/ARXIV.2101.08195

717    Chen, T., & Guestrin, C. (2016). *XGBoost: A Scalable Tree Boosting System*.

718         https://doi.org/10.48550/ARXIV.1603.02754

719    Chen, X., Huang, J., Han, Z., Gao, H., Liu, M., Li, Z., Liu, X., Li, Q., Qi, H., & Huang, Y.

720         (2020). The importance of short lag-time in the runoff forecasting model based

721         on long short-term memory. *Journal of Hydrology*, *589*, 125359.

722         https://doi.org/10.1016/j.jhydrol.2020.125359

723    Chen, Z., Badrinarayanan, V., Lee, C.-Y., & Rabinovich, A. (2017). *GradNorm: Gradient*

724         *Normalization for Adaptive Loss Balancing in Deep Multitask Networks*.

725         https://doi.org/10.48550/ARXIV.1711.02257

726    Choulga, M., Kourzeneva, E., Balsamo, G., Boussetta, S., & Wedi, N. (2019). Upgraded

727         global mapping information for earth system modelling: An application to

728    surface water depth at the ECMWF. *Hydrology and Earth System Sciences*,

729    *23*(10), 4051–4076. https://doi.org/10.5194/hess-23-4051-2019

730    Cui, Z., Qing, X., Chai, H., Yang, S., Zhu, Y., & Wang, F. (2021). Real-time rainfall-runoff

731    prediction using light gradient boosting machine coupled with singular spectrum

732    analysis. *Journal of Hydrology*, *603*, 127124.

733    https://doi.org/10.1016/j.jhydrol.2021.127124

734    Datta, P., & Faroughi, S. A. (2023). A multihead LSTM technique for prognostic prediction

735    of soil moisture. *Geoderma*, *433*, 116452.

736    https://doi.org/10.1016/j.geoderma.2023.116452

737    De Rosnay, P., Balsamo, G., Albergel, C., Muñoz-Sabater, J., & Isaksen, L. (2014).

738    Initialisation of Land Surface Variables for Numerical Weather Prediction.

739    *Surveys in Geophysics*, *35*(3), 607–621. https://doi.org/10.1007/s10712-012-

740    9207-x

741    ECMWF. (n.d.). Forecast User Guide. In *Anomaly Correlation Coefficient*. Retrieved 4

742    July 2024, from

743    https://confluence.ecmwf.int/display/FUG/Section+6.2.2+Anomaly+Correlation

744    +Coefficient

745    ECMWF. (2017). *IFS Documentation CY43R3 - Part IV: Physical processes*.

746    https://doi.org/10.21957/EFYK72KL

747    ECMWF. (2023). *IFS Documentation CY48R1 - Part IV: Physical Processes*.

748    https://doi.org/10.21957/02054F0FBF

749    Fer, I., Kelly, R., Moorcroft, P. R., Richardson, A. D., Cowdery, E. M., & Dietze, M. C.

750    (2018). Linking big models to big data: Efficient ecosystem model calibration

751        through Bayesian model emulation. *Biogeosciences*, *15*(19), 5801–5830.

752        https://doi.org/10.5194/bg-15-5801-2018

753   Girshick, R. (2015). *Fast R-CNN*. https://doi.org/10.48550/ARXIV.1504.08083

754   Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. The MIT Press.

755   Grinsztajn, L., Oyallon, E., & Varoquaux, G. (2022). *Why do tree-based models still*

756        *outperform deep learning on tabular data?* (Version 1). arXiv.

757        https://doi.org/10.48550/ARXIV.2207.08815

758   Guo, Y., Yu, X., Xu, Y.-P., Chen, H., Gu, H., & Xie, J. (2021). AI-based techniques for multi-

759        step streamflow forecasts: Application for multi-objective reservoir operation

760        optimization and performance assessment. *Hydrol. Earth Syst. Sci.*

761   Hatfield, S., Chantry, M., Dueben, P., Lopez, P., Geer, A., & Palmer, T. (2021). Building

762        Tangent-Linear and Adjoint Models for Data Assimilation With Neural Networks.

763        *Journal of Advances in Modeling Earth Systems*, *13*(9), e2021MS002521.

764        https://doi.org/10.1029/2021MS002521

765   Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J.,

766        Nicolas, J., Peubey, C., Radu, R., Schepers, D., Simmons, A., Soci, C., Abdalla,

767        S., Abellan, X., Balsamo, G., Bechtold, P., Biavati, G., Bidlot, J., Bonavita, M., …

768        Thépaut, J. (2020). The ERA5 global reanalysis. *Quarterly Journal of the Royal*

769        *Meteorological Society*, *146*(730), 1999–2049. https://doi.org/10.1002/qj.3803

770   Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural*

771        *Computation*, *9*(8), 1735–1780. https://doi.org/10.1162/neco.1997.9.8.1735

772   James, W., & Stein, C. (1992). Estimation with Quadratic Loss. In S. Kotz & N. L. Johnson

773        (Eds.), *Breakthroughs in Statistics* (pp. 443–460). Springer New York.

774        https://doi.org/10.1007/978-1-4612-0919-5_30

775    Kadra, A., Lindauer, M., Hutter, F., & Grabocka, J. (2021). *Well-tuned Simple Nets Excel*

776          *on Tabular Datasets* (arXiv:2106.11189). arXiv. http://arxiv.org/abs/2106.11189

777    Keisler, R. (2022). *Forecasting Global Weather with Graph Neural Networks*

778          (arXiv:2202.07575). arXiv. http://arxiv.org/abs/2202.07575

779    Kimpson, T., Choulga, M., Chantry, M., Balsamo, G., Boussetta, S., Dueben, P., &

780          Palmer, T. (2023). Deep learning for quality control of surface physiographic

781          fields using satellite Earth observations. *Hydrology and Earth System Sciences*,

782          *27*(24), 4661–4685. https://doi.org/10.5194/hess-27-4661-2023

783    Kingma, D. P., & Ba, J. (2017). *Adam: A Method for Stochastic Optimization*.

784    Kratzert, F., Herrnegger, M., Klotz, D., Hochreiter, S., & Klambauer, G. (2019).

785          *NeuralHydrology—Interpreting LSTMs in Hydrology*.

786          https://doi.org/10.48550/ARXIV.1903.07903

787    Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G. (2019).

788          Towards learning universal, regional, and local hydrological behaviors via

789          machine learning applied to large-sample datasets. *Hydrology and Earth System*

790          *Sciences*, *23*(12), 5089–5110. https://doi.org/10.5194/hess-23-5089-2019

791    Lam, R., Sanchez-Gonzalez, A., Willson, M., Wirnsberger, P., Fortunato, M., Alet, F.,

792          Ravuri, S., Ewalds, T., Eaton-Rosen, Z., Hu, W., Merose, A., Hoyer, S., Holland, G.,

793          Vinyals, O., Stott, J., Pritzel, A., Mohamed, S., & Battaglia, P. (2023). Learning

794          skillful medium-range global weather forecasting. *Science*, *382*(6677), 1416–

795          1421. https://doi.org/10.1126/science.adi2336

796    Lang, S., Alexe, M., Chantry, M., Dramsch, J., Pinault, F., Raoult, B., Clare, M. C. A.,

797          Lessig, C., Maier-Gerber, M., Magnusson, L., Bouallègue, Z. B., Nemesio, A. P.,

798          Dueben, P. D., Brown, A., Pappenberger, F., & Rabier, F. (2024). *AIFS - ECMWF's*

799        *data-driven forecasting system (arXiv:2406.01465). arXiv.*

800        http://arxiv.org/abs/2406.01465

801   Lees, T., Reece, S., Kratzert, F., Klotz, D., Gauch, M., De Bruijn, J., Kumar Sahu, R., Greve,

802        P., Slater, L., & Dadson, S. J. (2022). Hydrological concept formation inside long

803        short-term memory (LSTM) networks. *Hydrology and Earth System Sciences*,

804        *26*(12), 3079–3101. https://doi.org/10.5194/hess-26-3079-2022

805   Li, L., Carver, R., Lopez-Gomez, I., Sha, F., & Anderson, J. (2023). *SEEDS: Emulation of*

806        *Weather Forecast Ensembles with Diffusion Models.*

807        https://doi.org/10.48550/ARXIV.2306.14066

808   Machac, D., Reichert, P., & Albert, C. (2016). Emulation of dynamic simulators with

809        application to hydrology. *Journal of Computational Physics*, *313*, 352–366.

810        https://doi.org/10.1016/j.jcp.2016.02.046

811   Meyer, D., Grimmond, S., Dueben, P., Hogan, R., & Van Reeuwijk, M. (2022). Machine

812        Learning Emulation of Urban Land Surface Processes. *Journal of Advances in*

813        *Modeling Earth Systems*, *14*(3), e2021MS002744.

814        https://doi.org/10.1029/2021MS002744

815   Mironov, D., & Helmert, J. (n.d.). *Parameterization of Lakes in NWP and Climate Models.*

816   Muñoz-Sabater, J., Dutra, E., Agustí-Panareda, A., Albergel, C., Arduini, G., Balsamo, G.,

817        Boussetta, S., Choulga, M., Harrigan, S., Hersbach, H., Martens, B., Miralles, D.

818        G., Piles, M., Rodríguez-Fernández, N. J., Zsoter, E., Buontempo, C., & Thépaut,

819        J.-N. (2021). ERA5-Land: A state-of-the-art global reanalysis dataset for land

820        applications. *Earth System Science Data*, *13*(9), 4349–4383.

821        https://doi.org/10.5194/essd-13-4349-2021

822    Nath, S., Lejeune, Q., Beusch, L., Seneviratne, S. I., & Schleussner, C.-F. (2022).

823        MESMER-M: An Earth system model emulator for spatially resolved monthly

824        temperature. *Earth System Dynamics*, *13*(2), 851–877.

825        https://doi.org/10.5194/esd-13-851-2022

826    Nearing, G., Cohen, D., Dube, V., Gauch, M., Gilon, O., Harrigan, S., Hassidim, A., Klotz,

827        D., Kratzert, F., Metzger, A., Nevo, S., Pappenberger, F., Prudhomme, C., Shalev,

828        G., Shenzis, S., Tekalign, T. Y., Weitzner, D., & Matias, Y. (2024). Global prediction

829        of extreme floods in ungauged watersheds. *Nature*, *627*(8004), 559–563.

830        https://doi.org/10.1038/s41586-024-07145-1

831    Ozaki, Y., Tanigaki, Y., Watanabe, S., Nomura, M., & Onishi, M. (2022). Multiobjective

832        Tree-Structured Parzen Estimator. *Journal of Artificial Intelligence Research*, *73*,

833        1209–1250. https://doi.org/10.1613/jair.1.13188

834    Pappenberger, F., Ramos, M. H., Cloke, H. L., Wetterhall, F., Alfieri, L., Bogner, K.,

835        Mueller, A., & Salamon, P. (2015). How do I know if my forecasts are better?

836        Using benchmarks in hydrological ensemble prediction. *Journal of Hydrology*,

837        *522*, 697–713. https://doi.org/10.1016/j.jhydrol.2015.01.024

838    Popov, S., Morozov, S., & Babenko, A. (2019). *Neural Oblivious Decision Ensembles for*

839        *Deep Learning on Tabular Data* (Version 2). arXiv.

840        https://doi.org/10.48550/ARXIV.1909.06312

841    Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., &

842        Prabhat. (2019). Deep learning and process understanding for data-driven Earth

843        system science. *Nature*, *566*(7743), 195–204. https://doi.org/10.1038/s41586-

844        019-0912-1

845 Sener, O., & Koltun, V. (2018). *Multi-Task Learning as Multi-Objective Optimization*

846 *(Version 2).* arXiv. https://doi.org/10.48550/ARXIV.1810.04650

847 Shwartz-Ziv, R., & Armon, A. (2021). *Tabular Data: Deep Learning is Not All You Need*.

848 https://doi.org/10.48550/ARXIV.2106.03253

849 Thorpe, A., Bauer, P., Magnusson, L., & Richardson, D. (2013). *An evaluation of recent*

850 *performance of ECMWF's forecasts.* https://doi.org/10.21957/HI1EEKTR

851 Van Katwyk, P., Fox-Kemper, B., Seroussi, H., Nowicki, S., & Bergen, K. J. (2023). A

852 Variational LSTM Emulator of Sea Level Contribution From the Antarctic Ice

853 Sheet. *Journal of Advances in Modeling Earth Systems*, *15*(12), e2023MS003899.

854 https://doi.org/10.1029/2023MS003899

855 Viterbo, P. (2002). *Land_surface_processes*.

856 Wesselkamp, M., Moser, N., Kalweit, M., Boedecker, J., & Dormann, C. F. (2022).

857 *Process-guidance improves predictive performance of neural networks for*

858 *carbon turnover in ecosystems.* https://doi.org/10.48550/ARXIV.2209.14229

859 Zwart, J. A., Oliver, S. K., Watkins, W. D., Sadler, J. M., Appling, A. P., Corson-Dosch, H.

860 R., Jia, X., Kumar, V., & Read, J. S. (2023). Near-term forecasts of stream

861 temperature using deep learning and data assimilation in support of

862 management decisions. *JAWRA Journal of the American Water Resources*

863 *Association*, *59*(2), 317–337. https://doi.org/10.1111/1752-1688.13093

864