

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**

hat formatiert: Schriftart: 12 Pt.

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

hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.

13 4 BrainLinks-BrainTools, University of Freiburg, Germany

hat formatiert: Schriftart: (Standard) Times New Roman

14 5 World Meteorological Organization, Geneva, Switzerland

hat formatiert: Schriftart: 12 Pt.

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.**

- hat gelöscht: and
- hat gelöscht: , and coupled with an atmospheric model provide boundary and initial conditions
- hat formatiert: Schriftart: (Standard) Times New Roman
- hat gelöscht: being
- hat gelöscht: ,
- hat gelöscht: ed
- hat gelöscht: progress in
- hat gelöscht: in
- hat gelöscht: by simulating land surface processes, which are integral to forecasting water, carbon, and energy fluxes in coupled atmospheric models
- hat gelöscht: the
- hat gelöscht: The
- hat gelöscht: , yet reliable
- hat gelöscht: numerical
- hat gelöscht: .

58 1 Introduction

59
60 While forecasting of climate and weather system processes has long been a task for numerical
61 models, recent developments in deep learning have introduced competitive machine-learning
62 (ML) systems for numerical weather prediction (NWP) (Bi et al., 2022; Lam et al., 2023;
63 Lang et al., 2024). Land surface models (LSMs), even though being an integral part of
64 numerical weather prediction, have not yet caught the attention of the ML-community. LSMs
65 forecast water, carbon and energy fluxes and, in coupling with an atmospheric model, provide
66 the lower boundary and initial conditions (Boussetta et al., 2021; De Rosnay et al.,
67 2014). The parametrization of land surface states does not only affect predictability of earth
68 and climate systems on sub-seasonal scales (Muñoz-Sabater et al., 2021), but also the short-
69 and medium-range skill of NWP forecasts (De Rosnay et al., 2014). Beyond their online
70 integration with NWPs, offline versions of LSMs provide research tools for experiments on
71 the land surface (Boussetta et al., 2021), the diversity of which, however, are limited by
72 substantial computational resources requirements and often moderate runtime efficiencies
73 (Reichstein et al., 2019).
74 Emulators constitute statistical surrogates for numerical simulation models that, by
75 approximating the latter, aim for increasing computational efficiency (Machac et al., 2016).
76 While the construction of emulators can itself require substantial computational resources,
77 their subsequent evaluation usually runs orders of magnitude faster than the original
78 numerical model (Fer et al., 2018). For this reason, emulators have found application for
79 example in modular parametrization of online weather forecasting systems (Chantry et al.,
80 2021), in replacing the MCMC-sampling procedure in Bayesian calibration of ecosystem
81 models (Fer et al., 2018), or in generating forecast ensembles of atmospheric states for
82 uncertainty quantification (Li et al., 2023). Beyond their computational efficiency, surrogate
83 models with high parametric flexibility have the potential to correct process mis-specification
84 in a physical model when fine-tuned to observations (Wesselkamp et al., 2022).
85 Modelling approaches used for emulation range from low parametrized, auto-regressive
86 linear models to highly non-linear and flexible neural networks (Baker et al., 2022; Chantry
87 et al., 2021; Meyer et al., 2022; Nath et al., 2022). In the global land surface system M-
88 MESMER, a set of simple AR1 regression models is used to initialize the numerical LSM,
89 resulting in a modularized emulator (Nath et al., 2022). Numerical forecasts of gross primary
90 productivity and hydrological targets were successfully approximated by Gaussian processes
91 (Baker et al., 2022; Machac et al., 2016), the advantage of which is their direct quantification

hat formatiert: Schriftart: 12 Pt.

hat gelöscht: the

hat gelöscht: a

hat gelöscht: s

hat formatiert: Schriftart: (Standard) Times New Roman

hat gelöscht: , (Lang et al., 2024)

hat gelöscht: s,

hat gelöscht: [3], [4]

hat gelöscht: thus

hat formatiert: Schriftart: (Standard) Times New Roman

hat formatiert: Schriftart: (Standard) Times New Roman

hat gelöscht: whic

hat gelöscht: h

hat gelöscht: are

hat formatiert: Schriftart: (Standard) Times New Roman

hat gelöscht: the required

hat formatiert: Schriftart: (Standard) Times New Roman

hat gelöscht: at

hat formatiert: Schriftart: (Standard) Times New Roman

hat gelöscht: for

hat gelöscht:

hat gelöscht: themselves

hat gelöscht:

hat formatiert: Schriftart: (Standard) Times New Roman

hat formatiert: Schriftart: (Standard) Times New Roman

hat formatiert: Schriftart: (Standard) Times New Roman

hat gelöscht: forecast

hat formatiert: Schriftart: (Standard) Times New Roman

hat gelöscht: for

hat gelöscht: and

hat gelöscht: improve predictions towards

hat formatiert: Schriftart: (Standard) Times New Roman

hat formatiert: Schriftart: (Standard) Times New Roman

hat gelöscht: , (Baker et al., 2022), (Chantry et al., 2021), (Meyer et al., 2022)

hat formatiert: Schriftart: (Standard) Times New Roman

hat gelöscht: (Machac et al., 2016)

hat formatiert: Schriftart: (Standard) Times New Roman

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

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

136 Long-short term memory networks (LSTMs) have gained popularity in hydrological
 137 forecasting as rainfall-runoff models, for predicting stream flow temperature and also soil
 138 moisture (Bassi et al., 2024; Kratzert, Klotz, et al., 2019; Lees et al., 2022; Zwart et al.,
 139 2023). Research on the interpretability of LSTMs has found correlations between the model
 140 cell states and spatially or thematically similar hydrological units (Lees et al., 2022),
 141 suggesting the specific usefulness of LSTM for representing variables with dynamic storages
 142 and reservoirs (Kratzert, Hernegger, et al., 2019). As emulators, LSTMs have been shown
 143 useful for sea surface level projection in a variational manner with Monte Carlo dropout (Van
 144 Katwyk et al., 2023).

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

- hat gelöscht:
- hat gelöscht: variables from
- hat gelöscht:
- hat gelöscht: to
- hat gelöscht: (Meyer et al., 2022)
- hat formatiert: Schriftart: (Standard) Times New Roman

- hat formatiert: Schriftart: (Standard) Times New Roman
- hat gelöscht: [5]
- hat gelöscht: , the advantage of which for the online framework is the temporal consistency of prognostic state variables...
- hat formatiert: Schriftart: (Standard) Times New Roman

- hat gelöscht: estimation

- hat gelöscht: [e.g.
- hat formatiert: Schriftart: (Standard) Times New Roman
- hat gelöscht: , (Lees et al., 2022), (Zwart et al., 2023), (Bassi et al., 2024)].
- hat formatiert: Schriftart: (Standard) Times New Roman
- hat formatiert: Schriftart: (Standard) Times New Roman
- hat formatiert: Schriftart: (Standard) Times New Roman
- hat formatiert: Schriftart: (Standard) Times New Roman

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

hat gelöscht: T

hat gelöscht:

hat gelöscht: evaluation

hat gelöscht: is done

hat gelöscht: will encompass transfer learning and validation on observations.

2 Methods

hat formatiert: Schriftart: 12 Pt.

2.1 The Land Surface Model: ECLand

180 ECLand is a tiled ECMWF Scheme for surface exchanges over land that represents surface
 181 heterogeneity and incorporates land surface hydrology (Balsamo et al., 2011; ECMWF,
 182 2017). ECLand computes surface turbulent fluxes of heat, moisture and momentum, and skin
 183 temperature over different tiles (vegetation, bare soil, snow, interception and water) and then
 184 calculates an area-weighted average for the grid-box to couple with the atmosphere
 185 (Boussetta et al., 2021). For the overall accuracy of the model, accurate land surface
 186 parameterizations are essential (Kimpson et al., 2023) as they e.g. determine the sensible and
 187 latent heat fluxes, and provide the lower boundary conditions for enthalpy and moisture
 188 equations in the atmosphere (Viterbo, 2002). We emulate three prognostic state variables of
 189 ECLand that represent core land surface processes: soil water volume ($m^3 m^{-3}$) and soil
 190 temperature (K) at each three depth layers (each at 0 – 7 cm, 7 – 21 cm and 21 – 72 cm) and
 191 snow cover fraction (%), aggregated at the surface layer.

hat gelöscht: S

hat gelöscht: E

hat gelöscht: L

hat gelöscht: (ECLand)

hat gelöscht: (ECMWF, 2017)

hat gelöscht: (

hat gelöscht:)

hat gelöscht: the land surface parameterization

hat gelöscht: s

hat gelöscht: ,

hat gelöscht: , so below are some more details on these parametrisations.

2.2 Data sources

213 As training data base, global simulation and reanalysis time series from 2010 to 2022 were
214 compiled to *zarr* format at an aggregated 6-hourly temporal resolution. Simulations and
215 climate fields were generated from ECMWFs development cycle CY49R2, ECLand and forced
216 by ERA-5 meteorological reanalysis data (Hersbach et al., 2020).
217 There are three main sources of data used for creation of the data base: The first is a selection
218 of surface physiographic fields from ERA5 (Hersbach et al., 2020) and their updated versions
219 (Boussetta et al., 2021; Choulga et al., 2019; Muñoz-Sabater et al., 2021), used as static
220 model input features (X). The second is a selection of atmospheric and surface model fields
221 from ERA5, used as static and dynamic model input features (Y). The third are ECLand
222 simulations, constituting the model's dynamic prognostic state variables (z) and hence
223 emulator input and target features. A total of 41 static, seasonal and dynamical features were
224 used to create the emulators, see table 1 for an overview of input variables and details on the
225 surface physiographic and atmospheric fields below.

hat gelöscht: l

hat gelöscht: , (Boussetta et al., 2021), (Muñoz-Sabater et al., 2021)

hat gelöscht: is

hat gelöscht: results

hat gelöscht: model

226 227 2.2.1 Surface physiographic fields

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

hat formatiert: Schriftart: 12 Pt.

hat gelöscht: then

hat gelöscht: to

hat gelöscht: e

hat formatiert: Schriftart: 12 Pt.

242 243 244 2.2.2 ERA5

hat formatiert: Schriftart: 12 Pt.

245

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

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

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

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat gelöscht: m3 m-3

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat gelöscht: m2

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

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

278

279 2.3 Emulators

280

281 We compare a long-short term memory neural network (LSTM), extreme gradient boosting
 282 regression trees (XGB) and a feedforward neural network (that we here refer to as multilayer
 283 perceptron, MLP). To motivate this setup and pave the way for discussing effects of (hyper-
 284)parameter choices, a short overview of all approaches is given. All analyses were conducted
 285 in Python. XGB was developed in dmlc's XGBoost python package¹. The MLP and LSTM
 286 were developed in the PyTorch lightning framework for deep learning². Neural networks
 287 were trained with the Adam algorithm for stochastic optimization (Kingma & Ba, 2017).
 288 Model architectures and algorithmic hyperparameters were selected through **combined**

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat gelöscht: m2

hat gelöscht: s-2

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat gelöscht: m2s

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat gelöscht: m2s

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat gelöscht: the utility of

hat gelöscht: that of

hat gelöscht: that of

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

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

296 Bayesian hyperparameter optimization with the Optuna framework (Akiba et al., 2019) and
297 additional manual tuning. The Bayesian optimization minimizes the neural network
298 validation accuracy, specified here as mean absolute error (MAE), over a predefined search
299 space for free hyperparameters with the Tree-structured Parzen Estimator (Ozaki et al., 2022).
300 The resulting hyperparameter and architecture choices which were used for the different
301 approaches are listed in the Supplementary Material.

302

303 2.3.1 MLP

304

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

314

315 2.3.2 LSTM

316

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

$$323 \quad \mathbf{h}_t, \mathbf{c}_t = f(\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{c}_{t-1}, \boldsymbol{\theta})$$

$$324 \quad \mathbf{z}_t = \mathbf{A}\mathbf{h}_t + \mathbf{b}$$

325 where \mathbf{h}_{t-1} denotes the hidden state, i.e. output estimates from the previous time step, \mathbf{c}_{t-1}
326 the cell state from the previous time step, and $\boldsymbol{\theta}$ the time-invariant model weights. We stacked
327 multiple LSTM cells to an encoder-decoder model with transfer layers for hidden and cell
328 state initialization and for transfer to the context vector (see figure 1) (Nearing et al., 2024). A

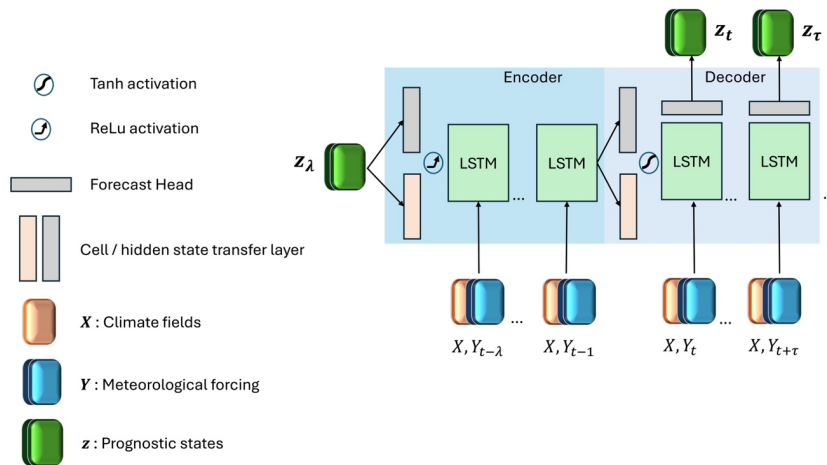
hat formatiert: Schriftart: (Standard) Times New Roman

hat formatiert: Schriftart: (Standard) Times New Roman

hat formatiert: Schriftart: (Standard) Times New Roman

hat formatiert: Schriftart: (Standard) Times New Roman

329 lookback l of the previous static and dynamic feature states are passed sequentially to the first
 330 LSTM cells in the encoder layer, while the l prognostic state variables \mathbf{z} initialize the hidden
 331 state \mathbf{h}_0 after a linear embedding. The output of the first LSTM layer cells become the input
 332 to the deeper LSTM layer cells and the last hidden state estimates are the final output from
 333 the encoder. Followed by a non-linear transformation with hyperbolic tangent activation, the
 334 hidden cell states are transformed into a weighted context vector \mathbf{s} . Together with the encoder
 335 the cell state (\mathbf{c}_t, \mathbf{s}) initializes the hidden and cell states of the decoder. The decoder LSTM
 336 cells take as input again static and dynamic features sequentially at lead times $t = 1, \dots, \tau$, but
 337 not the prognostic states variables. These are estimated from the sequential hidden states of
 338 the last LSTM layer cells, transformed to target size with a linear forecast head before
 339 prediction. LSTM predicts absolute state variables \mathbf{z}_t while being optimized on \mathbf{z}_t and $d\mathbf{z}_t$
 340 simultaneously, see section on optimization.



341
 342 *Figure 1: LSTM architecture. Blue shaded area indicates the encoder part, where the model*
 343 *is driven by a lookback λ of meteorological forcing and state variables. The light-blue shaded*
 344 *area indicates the decoder part that is initialized from the encoding to unroll LSTM forecasts*
 345 *from the initial time step t up to a flexibly long lead time of τ .*

346 **2.3.3 XGB**

347
 348 Extreme gradient boosting (XGB) is a regression tree ensemble method that uses an
 349 approximate algorithm for best split finding. It computes first and second order gradient
 350 statistics in the cost function, performing a similar to gradient descent optimization (T. Chen
 351 & Guestrin, 2016), where each new learner is trained on the residuals of the previous ones.

hat formatiert: Schriftart: (Standard) Times New Roman
 hat formatiert: Schriftart: (Standard) Times New Roman

hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.
 hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.
 hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.
 hat formatiert: Schriftart: 12 Pt.
 hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.
 hat formatiert: Schriftart: 12 Pt.
 hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.
 hat formatiert: Schriftart: 12 Pt.

352 Regularization and column sampling aim for preventing overfitting internally. XGB is known
353 to provide a powerful benchmark for time series forecasting and tabular data (T. Chen &
354 Guestrin, 2016; X. Chen et al., 2020; Shwartz-Ziv & Armon, 2021). Like the MLP, it is
355 trained to predict the increment $\Delta z_{t,i}$ of prognostic state variables, but only for a one-step
356 ahead prediction.

hat gelöscht: [

hat gelöscht: , (X. Chen et al., 2020)]

358 2.4 Experimental setup

359
360 We distinguish the experimental analysis into three parts that vary in the usage of the training
361 database: (1) model development, (2) model testing, and (3) global model transfer.

362 The models were developed and for the first time evaluated on a low state resolution
363 (ECMWF's TCO199 reduced gaussian grid, see section on data sources) and temporal subset
364 from the training data base, i.e. on a bounding box of 7715 grid cells over Europe with time
365 series of six years from 2016 to 2022. For details on the development data base, model
366 selection and model performances, see Supplementary Material S3.

367 The selected models were recreated on a high state resolution (TCO399) continental scale
368 European subset with 10 051 grid cells. Models were trained on five years 2015-2020 with
369 the year 2020 as validation split and evaluated on the year 2021 for the scores we report in
370 the main part. Note that for computation of forecast horizons, the two test years 2021 and
371 2022 were used, see details in section on forecast horizons. With this same data splitting
372 setup, the analysis was repeated in transferring the candidates to the low resolution (TCO199)
373 global data set with a total of 47892 grid cells. The low global resolution on one hand
374 allowed a systematic comparison of the three models, because high resolution training with
375 XGB was prohibited by the required working memory. On the other hand, this extrapolation
376 scenario created an unseen problem for the models that were selected on a continental and
377 high-resolution scale which is reflected in the resulting scores.

hat formatiert: Schriftart: 12 Pt.

379 2.5 Optimization

381 2.5.1 Loss functions

382

385 The basis of the loss function \mathcal{L} for the neural network optimization was PyTorch’s
 386 SmoothL1Loss³, a robust loss function that combines L1-norm and L2-norm and is less
 387 sensitive to outliers than pure L1-norm (Girshick, 2015). Based on a pre-defined threshold
 388 parameter β , smooth L1 transitions from L2-norm to L1-norm above the threshold.
 389 SmoothL1Loss \mathcal{L} is defined as

$$390 \quad \mathcal{L}(z, z) = 0.5(z - z)^2 \frac{1}{\beta} \text{ if } |z - z| < \beta \text{ and}$$

$$391 \quad \mathcal{L}(z, z) = |z - z| - 0.5 \beta \text{ otherwise,}$$

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

$$398 \quad \mathcal{L}_s(\hat{d}\mathbf{z}, d\mathbf{z}) = \sum_{t=1}^{\tau} \sum_{i=1}^p \mathcal{L}_t(\hat{d}\mathbf{z}_{t,i}, d\mathbf{z}_{t,i}),$$

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

400

$$401 \quad \mathcal{L}_r(\hat{d}\mathbf{z}, \mathbf{z}) = \frac{1}{\tau} \sum_{t=1}^{\tau} \sum_{i=1}^p \mathcal{L}_t(z_{t-1,i} + \hat{d}\mathbf{z}_{t,i}, z_{t,i})$$

402

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

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

411 **LSTM** The combined loss function for the LSTM is the sum of the incremental loss
 412 \mathcal{L}_s , where the $d\mathbf{z}_t$ were derived from \mathbf{z}_t after the forward pass, and the loss \mathcal{L} computed on

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

413 decoder estimates of prognostic states variables, a functionality that leverages the potential of
414 our LSTM structure.

415 **XGB** Trained only from one to the next time step, i.e. at a lead time of $\tau = 1$, the incremental
416 loss $\mathcal{L}_s = \mathcal{L}_r$. Without a SmoothL1Loss implementation provided in dmlc's XGBoost, we
417 trained XGB with both the Huber-Loss and the default L2-loss. The latter initially providing
418 better results, we chose the default L2-norm as loss function for XGB with the regularization
419 parameter $\lambda = 1$.

420

421 2.5.1 Normalization

422 As prognostic target variables are all lower bounded by zero, we tested both z-scoring and
423 max-scoring. The latter yielded no significant **improvement**; thus we show our results with z-
424 scored target variables. For neural network training but not for fitting XGB, static, dynamic
425 and prognostic state variables were all normalized with z-scoring towards the continental or
426 global mean z and unit standard deviation s_z as

$$427 z_{t,n} = \frac{z_{t,n} - z}{s_z}.$$

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

431

432 2.5.3 Spatial and temporal sampling

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

439

440 2.6 Evaluation

441

442 Three scores are used for model validation during the model development phase and in
443 validating architecture and hyperparameter selection, being the root mean squared error
444 (*RMSE*), the mean absolute error (*MAE*) and the anomaly correlation coefficient (*ACC*).

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

hat gelöscht: improvement,

hat gelöscht: backtransformed

448 against ECLand simulations directly with RMSE and MAE. Doing so, scores were
 449 aggregated over lead times, grid cells or both. The total RMSE was computed as

$$450 \quad \text{RMSE} = \sqrt{\frac{\sum_{\tau,p} (z - \hat{z})^2}{n}},$$

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

$$453 \quad \text{MAE} = \frac{\sum_{\tau,p} |z - \hat{z}|}{n},$$

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

$$461 \quad \text{ACC}(t) = \frac{\overline{(z - c)(\hat{z} - c)}}{\sqrt{\overline{(z - c)^2} \overline{(\hat{z} - c)^2}}}$$

462 at each $t = 1, \dots, \tau$ where the overbar denotes averaging over grid cells $p = i, \dots, n$. This way,
 463 the nominator represents the average spatial covariance of emulator and numerical forecasts
 464 with climatology as expected sample mean. Hence, it indicates the mean squared skill error
 465 towards climatology, and the denominator indicates its variability. The aggregated scores that
 466 are shown in tables 3-5 represent the temporally arithmetic mean of ACC(t). ACC is bounded
 467 between 1 and -1, and an ACC of 1 indicates perfect representation of forecast error
 468 variability, an ACC of 0.5 indicates a similar forecast error to that of the climatology, an ACC
 469 of 0 indicates that forecast error variability dominates and the forecast has no value and an
 470 ACC approaching -1 indicates that the forecast has been very unreliable (ECMWF, n.d.).
 471 ACC is undefined when the denominator is zero. This is the case either when mean squared
 472 emulator or ECLand anomaly, or both are zero because forecast and climatology perfectly
 473 align, or because they cancel out at summation to the mean.

474
 475 **2.6.1 Forecast horizons**
 476

hat gelöscht: now

hat gelöscht:

479 Forecast horizons of the emulators are defined by the decomposition of the RMSE
480 (Bengtsson et al., 2008) into the emulator's variability around climatology (i.e. anomaly),
481 ECLand's variability around climatology and the covariance of both. The horizon is the point
482 in time at which the forecast error reaches saturation level, that is when the covariance of
483 emulator and ECLand anomalies approaches zero, as does the ACC.
484 We analysed predictive ability and predictability by computing the ACC for all lead times
485 from 6 hours to approx. one year, i.e. lead times $t = 1, \dots, \tau$, τ being 1350. As this confounds
486 the seasonality with the lead time, we compute these for every starting point of the prediction,
487 requiring two test years (2021 and 2022).
488 Forecast horizons based on the emulators' skill in standardized anomaly towards persistence
489 were equivalently computed but with persistence as a benchmark for shorter time scales, this
490 was only done for three months, from January to March 2021.
491 The analysis was conducted on two exemplary regions in northern and southern Europe that
492 represent very different conditions orography and in prognostic land surface states,
493 specifically in snow cover. For details on the regions and on the horizons computed with
494 standardized anomaly skill, see Appendices A1 and A4 respectively.

496 **3 Results**

497
498 The improvement in evaluation runtimes achieved by emulators toward the numerical
499 ECLand were significant. Iterating the forecast over a full test year at 30 km spatial
500 resolution, XGB evaluates in 5.4 minutes, LSTM in 3.09 minutes and MLP in 0.05 minutes
501 (i.e. 3.2 seconds) on average. In contrast, ECLand integration over a full test year on 16
502 CPUs at 30 km spatial resolution takes approximately 240 minutes (i.e. four hours). The slow
503 runtime of the LSTM compared to the MLP emulator is caused by a spatial chunking
504 procedure that was not optimise for this work, but could be improved in the future.

506 **3.1 Aggregated performances**

507
508 **Europe.** All emulators approximated the numerical LSM with high average total accuracies
509 (all RMSEs < 1.58 and MAEs < 0.84) and confident correlations (all ACC > 0.72) (see table
510 2 and figure 2). The LSTM emulator achieved the best results across all total average scores
511 on the European scale. It decreased the total average MAE by ~25% towards XGB and by
512 ~37% towards the MLP and the total average RMSE by ~42% towards XGB and ~38%

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: Nicht Fett

hat formatiert: Schriftart: 12 Pt.

513 towards the MLP. In total average ACC, the LSTM scored 20% higher than the MLP and
514 15% than XGB, also being the only emulator that achieved an ACC > 0.9. While the MLP
515 outperforms XGB in total average RMSE by ~5%, XGB scores better than the MLP in MAE
516 by ~27%.

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

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

530 In snow cover approximation, the LSTM emulator enhances accuracies by over ~50% in
531 MAE towards both the XGB and the MLP emulator and scores highest in anomaly
532 representation with an ACC of ~0.87 compared to an ACC of ~0.66 for the MLP and only
533 ~0.74 for the XGB (see table 5).

534 **Globe.** Score ranking on the global scale varies strongly from the continental scale (see table
535 2). In total average accuracies, the MLP outperforms XGB by over 30% and LSTM by up
536 ~25% in RMSE and improves MAE more than 15% towards both. In anomaly correlation
537 however it scores last, whereas XGB achieves the highest total average of over 0.75.

538 Consistent with scores on the continental scale is XGBs high performance in soil temperature
539 (see table 3). It significantly outperforms the LSTM by ~60% (K) in RMSE and nearly up to
540 75% (K) in MAE in all layers and the MLP by up to 50% (K) in MAE at the top layer.

541 Anomaly persistence for all models degrade visibly towards the lower soil layers, while that
542 of the LSTM most relative to MLP and XGB. Like on the continental scale, XGB also
543 outperforms the other candidates in soil temperature forecasts in all but the medium layer,
544 where the MLP gets higher scores in MAE and RMSE but not in ACC (see table 4). LSTM
545 doesn't stand out with any scores on the global scale.

546

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat gelöscht:

hat gelöscht: Similar to

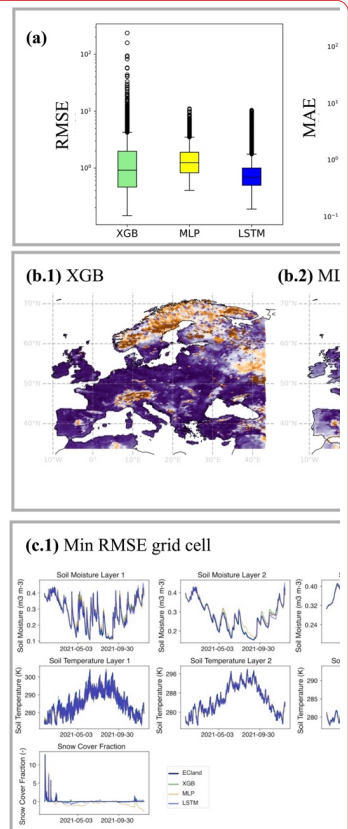
549 **3.2 Spatial and temporal performances**

550

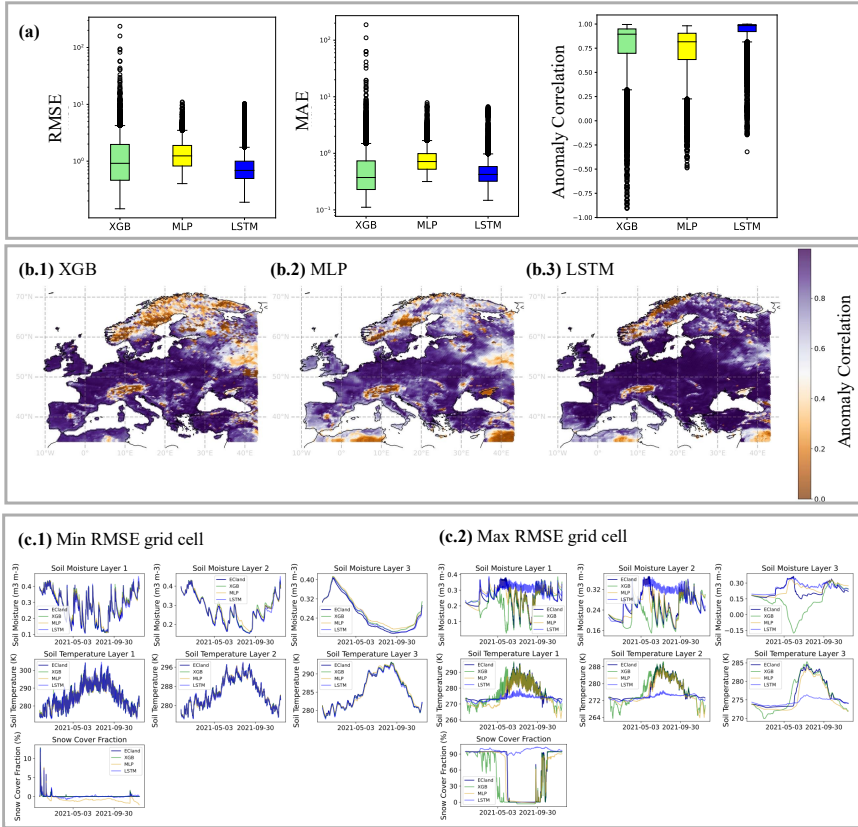
551 **Europe.** When summarizing temporally aggregated scores as boxplots to a total distribution
552 over space (see figure 2, A), the long tails of XGB scores become visible, whereas the MLP
553 indicates most robustness. This is reflected in the geographic distribution of scores at the
554 example of ACC (see figure 2, bottom), where the area of low anomaly correlation is largest
555 for XGB, ranging over nearly all northern Scandinavia, while MLP and LSTM have smaller
556 and more segregated areas of clearly low anomaly correlation. The LSTM shows a
557 homogenously high ACCs over most of central Europe but the Alps, while also seems to be
558 challenged in areas of relative to the central Europe extreme weather conditions at the
559 Norwegian and Spanish coasts.

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

hat gelöscht: Similar to



hat gelöscht:



569
 570 *Figure 2: a: Total aggregated distributions of (log) scores averaged over lead times, i.e.*
 571 *displaying the variation among grid cells. b: The distribution of the anomaly correlation in*
 572 *space on the European subset (b.1: XGB, b.2: MLP, b.3: LSTM). c: Model forecasts over test*
 573 *year 2021 for grid cell with minimum and maximum RMSE values (LSTM).*

574
 575 *Table 2: Emulator total average scores (unitless), aggregated over variables, time and space*
 576 *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

hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.

hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.

hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.

hat formatiert: Schriftart: (Standard) Times New Roman

hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

577 *Table 3: Emulator average scores (RMSE, MAE in $m^3 m^{-3}$) on soil water volume forecasts*
 578 *for the European subset, aggregated over space and time from the European and Global*
 579 *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

580
 581 *Table 4: Emulators' mean scores (RMSE, MAE in K) on soil temperature forecasts for the*
 582 *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	

583
 584 *Table 5: Emulators' mean scores (RMSE, MAE in %) on snow cover forecasts for the*
 585 *European subset, aggregated over space and time.*

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

hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.

hat gelöscht:

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

Snow top	XGB	8.219	9.906	3.099	5.196	0.746	0.707
cover	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

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

hat formatiert: Schriftart: 12 Pt.

587

588 3.3 Forecast horizons

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

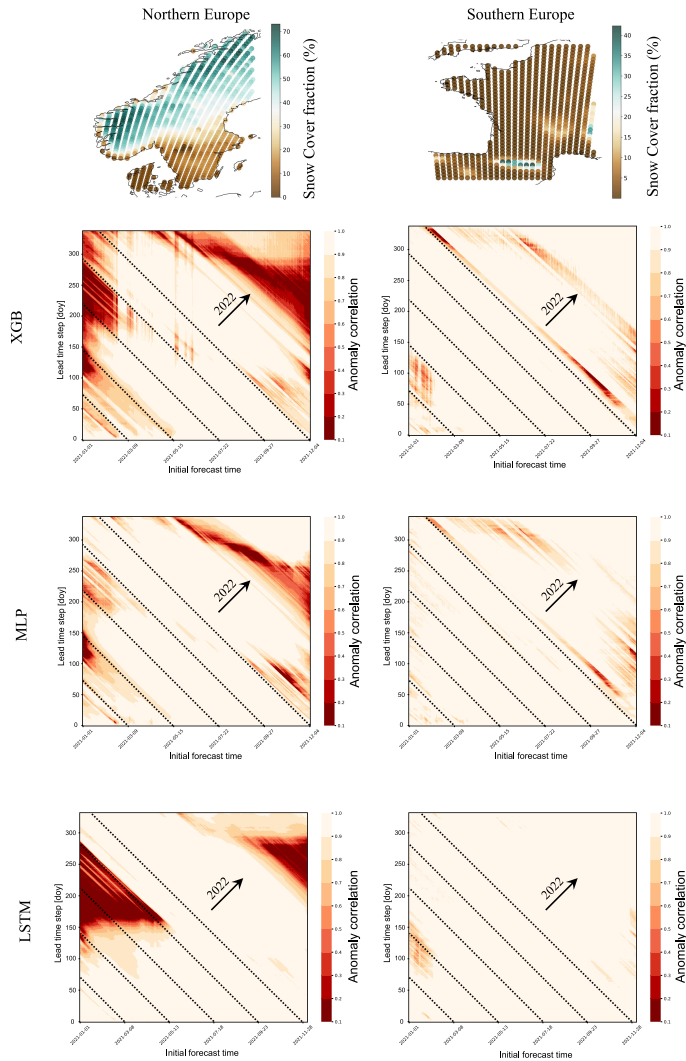
596 In the North, predictive skill depended on an interaction of how far ahead a prediction was
597 made (the lead time) and the day of year to which the prediction was made. In the best case,
598 the LSTM, summer predictions were poor (light patches in figure 3 heat maps), but only
599 when initialised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

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

hat gelöscht: d

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

620



621

622 *Figure 3: Top row, European subregions for computations of forecast skill horizons and their*
 623 *yearly average snow cover fraction (%), predicted by ECLand. Rows 2-4, Emulator forecast*
 624 *skill horizons in the subregions, aggregated over prognostic state variables, computed with*

hat gelöscht:

hat formatiert: Schriftart: (Standard) Times New Roman

hat formatiert ... [1]

hat gelöscht: in ... n the two ... European ... ubregions, aggregated over prognostic state variables, s. Scores at ... [2]

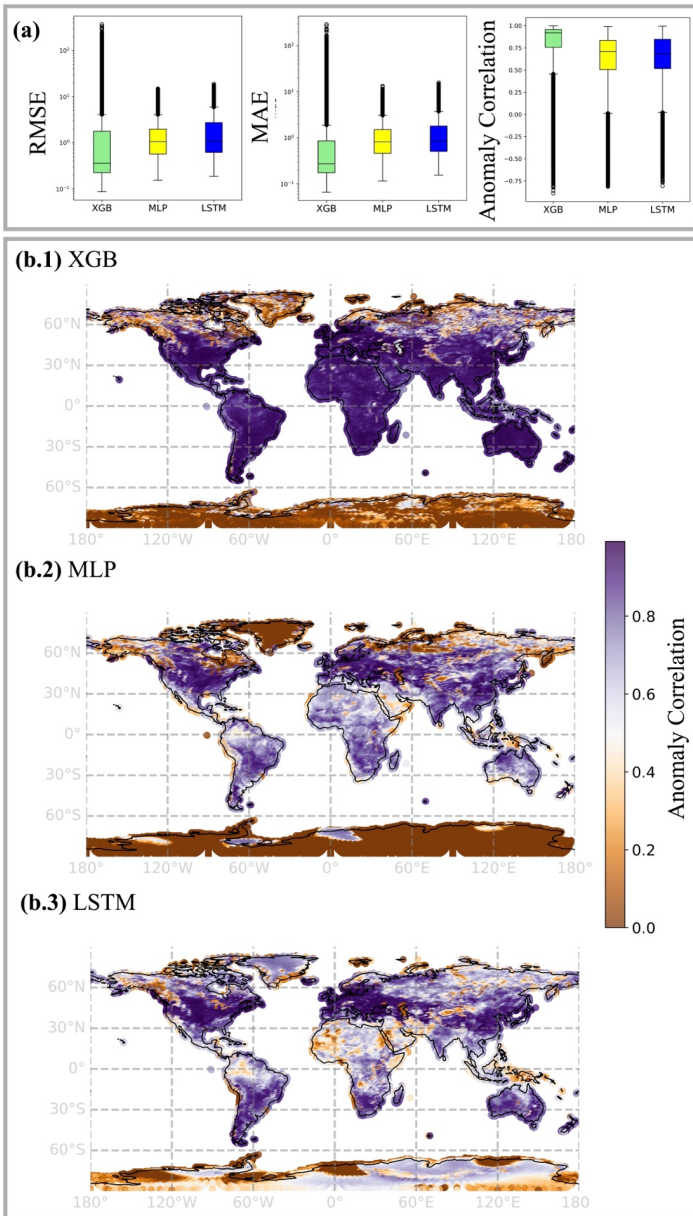
633 the anomaly correlation coefficient (ACC) at 6-hourly lead times (y-axis) over approx. one
634 year, displayed as a function of the initial forecast time (x-axis). *The horizon is the time at*
635 *which the forecast has no value at all, i.e. when ACC is 0 (or below 10%). The diagonal*
636 *dashed lines indicate the day of the test year 2021 as labelled on the x-axis, the arrows*
637 *indicate where forecasts reach the second test year 2022.*

638

hat gelöscht: As

hat gelöscht: we define

hat formatiert: Schriftart: 12 Pt.



641
642
643
644

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

hat formatiert: Schriftart: (Standard) Times New Roman

hat formatiert: Schriftart: (Standard) Times New Roman

hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.

hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.

hat formatiert: Schriftart: (Standard) Times New Roman, 12 Pt.

hat formatiert: Schriftart: 12 Pt.

645 **4 Discussion**

646

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

658

659 **4.1 Approximation of prognostic land surface states**

660

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

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

hat gelöscht: emulators

hat gelöscht: models

hat gelöscht: With

hat gelöscht: want to

hat gelöscht:

hat gelöscht: in the

685 Inversely to expectations and preceding experiments, on the European data set relative to the
686 two other models the LSTM scored better in the upper layer soil temperatures than in
687 forecasting soil water volume and decreased in scores towards lower layers with slower
688 processes. For training on observations, the decreasing LSTM predictive accuracy for soil
689 moisture with lead time is discussed (Datta & Faroughi, 2023), but reasons arising from the
690 engineering side remain unclear. In an exemplary case of a single-objective, deterministic
691 streamflow forecast, a decrease in recurrent neural network performance has been related
692 with an increasing coefficient of variation (Y. Guo et al., 2021). In our European subregions,
693 the signal-to-noise ratio of the prognostic state variables (computed as the averaged ratio of
694 mean and standard deviation) is up to ten times higher in soil temperature than in soil water
695 volume states (see Supplementary Material, S2.1). While a small signal of the latter may
696 induce instability in scores, it does not explain the decreasing performance towards deeper
697 soil layers with slow processes, where we expected an advantage of the long-term memory.
698 Stein's paradox tells us that joint optimization may lead to better results if the target is multi-
699 objective, but not if we are interested in single targets (James & Stein, 1992; Sener & Koltun,
700 2018). While from a process perspective multi-objective scores are less meaningful than
701 single ones, this is what we opted for due to efficiency. The unweighted linear loss
702 combination might be suboptimal in finding effective parameters across all prognostic state
703 variables (Z. Chen et al., 2017; Sener & Koltun, 2018), yet being strongly correlated, we
704 deemed their manual weighting inappropriate. An alternative to this provides adaptive loss
705 weighting with gradient normalisation (Z. Chen et al., 2017).

hat gelöscht: (Sener & Koltun, 2018)

hat gelöscht: (Sener & Koltun, 2018)

707 4.2 Evaluation in time and space

708
709 We used aggregated MAE and RMSE accuracies as a first assessment tool to conduct model
710 comparison, but score aggregation hides model specific spatio-temporal residual patterns.
711 Further, both scores are variance dependent, favouring low variability in model forecasts
712 even though this may not be representative of the system dynamic (Thorpe et al., 2013).
713 Assessing the forecast skill over time as the relative proximity to a subjectively chosen
714 benchmark helps disentangling areas of strengths and weaknesses in forecasting with the
715 emulators (Pappenberger et al., 2015). The naïve 6-hourly climatology as benchmark
716 highlights periods where emulators long-range forecasts on the test year are externally limited
717 by seasonality, i.e. system predictability, and where they are internally limited by model error,
718 i.e. the model's predictive ability. Applying this strategy in two exemplary European

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

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

hat gelöscht: we envision

hat gelöscht: for

hat gelöscht: follow-up

735 4.3 Emulation with memory mechanisms

736
737 Without much tuning, XGB challenges both LSTM and MLP for nearly all variables (see
738 tables 2-4). In training on observations for daily short-term and real-time rainfall-runoff
739 prediction, XGB and LightXGB were shown before to equally performed as, or outperformed
740 LSTMs (X. Chen et al., 2020; Cui et al., 2021). Nevertheless, models with memory
741 mechanism such as the encoder-decoder LSTM remain a promising approach for land surface
742 forecasting regarding their differentiability (Hatfield et al., 2021), their flexible extension of
743 lead times, for exploring the effect of long-term dependencies or for inference from the
744 context vector that may help identifying the process relevant climate fields (Lees et al.,
745 2022).

hat gelöscht: (Cui et al., 2021)

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

hat gelöscht: In our

hat gelöscht:

hat gelöscht: , we

hat gelöscht: our

753 Comparing average predictive accuracies across different training lead times indicates that
754 training at longer lead times may enhance short-term accuracy of the LSTM at the cost of

763 training runtime (see Supplementary Material, S2). A superficial exploration of encoder
764 length indicates no visible improvement on target accuracies if not a positive tendency
765 towards shorter sequences. This needs an extended analysis for understanding, yet without a
766 significant improvement by increased sequence length, GRU cells might provide a simplified
767 and less parameterized alternative to LSTM cells. They were found to perform equally well
768 on streamflow forecast performance before, while reaching higher operational speed (Y. Guo
769 et al., 2021).

770

771 **4.4 Emulators in application**

772

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

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

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

hat gelöscht: (Grinsztajn et al., 2022)

798 differentiable, limiting its application in coupled data assimilation (Hatfield et al., 2021) even
799 though research on differentiable trees is ongoing (Popov et al., 2019).

800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830

4.5 Experimentation with Emulators

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

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

- hat formatiert: Schriftart: 12 Pt.
- hat formatiert: Schriftart: (Standard) Times New Roman
- hat formatiert: Schriftart: (Standard) Times New Roman
- hat formatiert: Schriftart: (Standard) Times New Roman
- hat formatiert: Schriftart: (Standard) Times New Roman
- hat formatiert: Schriftart: (Standard) Times New Roman
- hat formatiert: Schriftart: (Standard) Times New Roman
- hat formatiert: Schriftart: (Standard) Times New Roman

Formatiert: Zeilenabstand: 1,5 Zeilen

- hat gelöscht: (Beucler et al., 2021; Brenowitz et al., 2020)
- hat formatiert: Schriftart: (Standard) Times New Roman
- hat formatiert: Schriftart: 12 Pt.

832 **5 Conclusion**

833

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

hat gelöscht: In conclusion

847

848 **Code and data availability**

849 Code for this analysis is published at [https://github.com/MWesselkamp/land-surface-](https://github.com/MWesselkamp/land-surface-emulation)
850 [emulation](https://github.com/MWesselkamp/land-surface-emulation). Training data is published at [10.21957/n17n-6a68](https://doi.org/10.21957/n17n-6a68) (Tco199) and [10.21957/pcf3-](https://doi.org/10.21957/pcf3-ah06)
851 [ah06](https://doi.org/10.21957/pcf3-ah06) (Tco399).

hat gelöscht: By developing the emulators for ECMWF's numerical land surface scheme ECLand, we path the way towards a physics-informed ML-based land surface model that on the long run can be parametrized with observations and provide a pretrained model suite to improve land surface forecasts.

hat formatiert: Schriftart: 12 Pt.

hat gelöscht: can be found

hat gelöscht: here:

hat formatiert: Schriftart: (Standard) Times New Roman

hat gelöscht: D

hat gelöscht: is available on request.

hat formatiert: Schriftart: 12 Pt.

852 **Author contribution**

853 MW, MCha, EP, FP and GB conceived the study. MW and EP conducted the analysis. MW,
854 MCha, MK, EP discussed and took technical decisions. SB advised on process decisions.
855 MW, MCho and FP wrote the manuscript. MW, MCha, EP, MCho, SB, MK, CFD, FP
856 reviewed the analysis and/or manuscript.

857 **Competing interest**

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

hat formatiert: Schriftart: 12 Pt.

859 **Acknowledgements**

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

hat formatiert: Schriftart: 12 Pt.

876 [reflect those of the European Union or the Commission](#). ChatGPT version 4.0 was used for
877 coding support.

878 ▲ hat formatiert: Schriftart: 12 Pt.

879 **References**

- 880
- 881 Akiba, T., Sano, S., Yanase, T., Ohta, T., & Koyama, M. (2019). Optuna: A Next-generation
882 Hyperparameter Optimization Framework. *Proceedings of the 25th ACM SIGKDD
883 International Conference on Knowledge Discovery & Data Mining*, 2623–2631.
884 <https://doi.org/10.1145/3292500.3330701>
- 885 Baker, E., Harper, A. B., Williamson, D., & Challenor, P. (2022). Emulation of high-
886 resolution land surface models using sparse Gaussian processes with
887 application to JULES. *Geoscientific Model Development*, 15(5), 1913–1929.
888 <https://doi.org/10.5194/gmd-15-1913-2022>
- 889 Balsamo, G., Boussetta, S., Dutra, E., Beljaars, A., & Viterbo, P. (2011). *Evolution of land-
890 surface processes in the IFS*. 127.
- 891 Bassi, A., Höge, M., Mira, A., Fenicia, F., & Albert, C. (2024). *Learning Landscape
892 Features from Streamflow with Autoencoders*. [https://doi.org/10.5194/hess-
893 2024-47](https://doi.org/10.5194/hess-2024-47)
- 894 Bengtsson, L. K., Magnusson, L., & Källén, E. (2008). Independent Estimations of the
895 Asymptotic Variability in an Ensemble Forecast System. *Monthly Weather
896 Review*, 136(11), 4105–4112. <https://doi.org/10.1175/2008MWR2526.1>
- 897 Beucler, T., Pritchard, M., Rasp, S., Ott, J., Baldi, P., & Gentine, P. (2021). Enforcing
898 Analytic Constraints in Neural Networks Emulating Physical Systems. *Physical
899 Review Letters*, 126(9), 098302. <https://doi.org/10.1103/PhysRevLett.126.098302>

900 Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X., & Tian, Q. (2022). *Pangu-Weather: A 3D High-*
901 *Resolution Model for Fast and Accurate Global Weather Forecast.*
902 <https://doi.org/10.48550/ARXIV.2211.02556>

903 Boussetta, S., Balsamo, G., Arduini, G., Dutra, E., McNorton, J., Choulga, M., Agustí-
904 Panareda, A., Beljaars, A., Wedi, N., Muñoz-Sabater, J., De Rosnay, P., Sandu, I.,
905 Hadade, I., Carver, G., Mazzetti, C., Prudhomme, C., Yamazaki, D., & Zsoter, E.
906 (2021). ECLand: The ECMWF land surface modelling system. *Atmosphere*, 12(6),
907 723. <https://doi.org/10.3390/atmos12060723>

908 Bouthillier, X., Delaunay, P., Bronzi, M., Trofimov, A., Nichyporuk, B., Szeto, J., Sepah, N.,
909 Raff, E., Madan, K., Voleti, V., Kahou, S. E., Michalski, V., Serdyuk, D., Arbel, T.,
910 Pal, C., Varoquaux, G., & Vincent, P. (2021). *Accounting for Variance in Machine*
911 *Learning Benchmarks* (arXiv:2103.03098). arXiv. <http://arxiv.org/abs/2103.03098>

912 Brenowitz, N. D., Henn, B., McGibbon, J., Clark, S. K., Kwa, A., Perkins, W. A., Watt-
913 Meyer, O., & Bretherton, C. S. (2020). *Machine Learning Climate Model*
914 *Dynamics: Offline versus Online Performance* (Version 1). arXiv.
915 <https://doi.org/10.48550/ARXIV.2011.03081>

916 Chantry, M., Hatfield, S., Duben, P., Polichtchouk, I., & Palmer, T. (2021). *Machine*
917 *learning emulation of gravity wave drag in numerical weather forecasting.*
918 <https://doi.org/10.48550/ARXIV.2101.08195>

919 Chen, T., & Guestrin, C. (2016). *XGBoost: A Scalable Tree Boosting System.*
920 <https://doi.org/10.48550/ARXIV.1603.02754>

921 Chen, X., Huang, J., Han, Z., Gao, H., Liu, M., Li, Z., Liu, X., Li, Q., Qi, H., & Huang, Y.
922 (2020). The importance of short lag-time in the runoff forecasting model based

923 on long short-term memory. *Journal of Hydrology*, 589, 125359.
924 <https://doi.org/10.1016/j.jhydrol.2020.125359>

925 Chen, Z., Badrinarayanan, V., Lee, C.-Y., & Rabinovich, A. (2017). *GradNorm: Gradient*
926 *Normalization for Adaptive Loss Balancing in Deep Multitask Networks*.
927 <https://doi.org/10.48550/ARXIV.1711.02257>

928 Choulga, M., Kourzeneva, E., Balsamo, G., Boussetta, S., & Wedi, N. (2019). Upgraded
929 global mapping information for earth system modelling: An application to
930 surface water depth at the ECMWF. *Hydrology and Earth System Sciences*,
931 23(10), 4051–4076. <https://doi.org/10.5194/hess-23-4051-2019>

932 Cui, Z., Qing, X., Chai, H., Yang, S., Zhu, Y., & Wang, F. (2021). Real-time rainfall-runoff
933 prediction using light gradient boosting machine coupled with singular spectrum
934 analysis. *Journal of Hydrology*, 603, 127124.
935 <https://doi.org/10.1016/j.jhydrol.2021.127124>

936 Datta, P., & Faroughi, S. A. (2023). A multihead LSTM technique for prognostic prediction
937 of soil moisture. *Geoderma*, 433, 116452.
938 <https://doi.org/10.1016/j.geoderma.2023.116452>

939 De Rosnay, P., Balsamo, G., Albergel, C., Muñoz-Sabater, J., & Isaksen, L. (2014).
940 Initialisation of Land Surface Variables for Numerical Weather Prediction.
941 *Surveys in Geophysics*, 35(3), 607–621. [https://doi.org/10.1007/s10712-012-](https://doi.org/10.1007/s10712-012-9207-x)
942 [9207-x](https://doi.org/10.1007/s10712-012-9207-x)

943 Dunkl, I., Spring, A., Friedlingstein, P., & Brovkin, V. (2021). Process-based analysis of
944 terrestrial carbon flux predictability. *Earth System Dynamics*, 12(4), 1413–1426.
945 <https://doi.org/10.5194/esd-12-1413-2021>

946 ECMWF. (n.d.). Forecast User Guide. In *Anomaly Correlation Coefficient*. Retrieved 4
947 July 2024, from
948 <https://confluence.ecmwf.int/display/FUG/Section+6.2.2+Anomaly+Correlation>
949 [+Coefficient](https://confluence.ecmwf.int/display/FUG/Section+6.2.2+Anomaly+Correlation)

950 ECMWF. (2017). *IFS Documentation CY43R3 - Part IV: Physical processes*.
951 <https://doi.org/10.21957/EFYK72KL>

952 ECMWF. (2023). *IFS Documentation CY48R1 - Part IV: Physical Processes*.
953 <https://doi.org/10.21957/02054F0FBF>

954 Fer, I., Kelly, R., Moorcroft, P. R., Richardson, A. D., Cowdery, E. M., & Dietze, M. C.
955 (2018). Linking big models to big data: Efficient ecosystem model calibration
956 through Bayesian model emulation. *Biogeosciences*, 15(19), 5801–5830.
957 <https://doi.org/10.5194/bg-15-5801-2018>

958 Gelbrecht, M., White, A., Bathiany, S., & Boers, N. (2023). Differentiable programming
959 for Earth system modeling. *Geoscientific Model Development*, 16(11), 3123–
960 3135. <https://doi.org/10.5194/gmd-16-3123-2023>

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

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

963 Grinsztajn, L., Oyallon, E., & Varoquaux, G. (2022). *Why do tree-based models still*
964 *outperform deep learning on tabular data?* (Version 1). arXiv.
965 <https://doi.org/10.48550/ARXIV.2207.08815>

966 Gu, M., Wang, X., & Berger, J. O. (2017). *Robust Gaussian Stochastic Process Emulation*
967 (Version 1). arXiv. <https://doi.org/10.48550/ARXIV.1708.04738>

968 Guo, Y., Yu, X., Xu, Y.-P., Chen, H., Gu, H., & Xie, J. (2021). AI-based techniques for multi-
969 step streamflow forecasts: Application for multi-objective reservoir operation
970 optimization and performance assessment. *Hydrol. Earth Syst. Sci.*

971 Guo, Z., Dirmeyer, P. A., & DelSole, T. (2011). Land surface impacts on subseasonal and
972 seasonal predictability: LAND IMPACTS SUBSEASONAL PREDICTABILITY.
973 *Geophysical Research Letters*, 38(24), n/a-n/a.
974 <https://doi.org/10.1029/2011GL049945>

975 Hatfield, S., Chantry, M., Dueben, P., Lopez, P., Geer, A., & Palmer, T. (2021). Building
976 Tangent-Linear and Adjoint Models for Data Assimilation With Neural Networks.
977 *Journal of Advances in Modeling Earth Systems*, 13(9), e2021MS002521.
978 <https://doi.org/10.1029/2021MS002521>

979 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J.,
980 Nicolas, J., Peubey, C., Radu, R., Schepers, D., Simmons, A., Soci, C., Abdalla,
981 S., Abellan, X., Balsamo, G., Bechtold, P., Biavati, G., Bidlot, J., Bonavita, M., ...
982 Thépaut, J. (2020). The ERA5 global reanalysis. *Quarterly Journal of the Royal*
983 *Meteorological Society*, 146(730), 1999–2049. <https://doi.org/10.1002/qj.3803>

984 Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural*
985 *Computation*, 9(8), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>

986 Irrgang, C., Boers, N., Sonnewald, M., Barnes, E. A., Kadow, C., Staneva, J., & Saynisch-
987 Wagner, J. (2021). Towards neural Earth system modelling by integrating artificial
988 intelligence in Earth system science. *Nature Machine Intelligence*, 3(8), 667–674.
989 <https://doi.org/10.1038/s42256-021-00374-3>

990 James, W., & Stein, C. (1992). Estimation with Quadratic Loss. In S. Kotz & N. L. Johnson
991 (Eds.), *Breakthroughs in Statistics* (pp. 443–460). Springer New York.
992 https://doi.org/10.1007/978-1-4612-0919-5_30

993 Kadra, A., Lindauer, M., Hutter, F., & Grabocka, J. (2021). *Well-tuned Simple Nets Excel*
994 *on Tabular Datasets* (arXiv:2106.11189). arXiv. <http://arxiv.org/abs/2106.11189>

995 Keisler, R. (2022). *Forecasting Global Weather with Graph Neural Networks*
996 (arXiv:2202.07575). arXiv. <http://arxiv.org/abs/2202.07575>

997 Kimpson, T., Choulga, M., Chantry, M., Balsamo, G., Boussetta, S., Dueben, P., &
998 Palmer, T. (2023). Deep learning for quality control of surface physiographic
999 fields using satellite Earth observations. *Hydrology and Earth System Sciences*,
1000 27(24), 4661–4685. <https://doi.org/10.5194/hess-27-4661-2023>

1001 Kingma, D. P., & Ba, J. (2017). *Adam: A Method for Stochastic Optimization*.

1002 Kratzert, F., Herrnegger, M., Klotz, D., Hochreiter, S., & Klambauer, G. (2019).
1003 *NeuralHydrology—Interpreting LSTMs in Hydrology*.
1004 <https://doi.org/10.48550/ARXIV.1903.07903>

1005 Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G. (2019).
1006 Towards learning universal, regional, and local hydrological behaviors via
1007 machine learning applied to large-sample datasets. *Hydrology and Earth System*
1008 *Sciences*, 23(12), 5089–5110. <https://doi.org/10.5194/hess-23-5089-2019>

1009 Lam, R., Sanchez-Gonzalez, A., Willson, M., Wirnsberger, P., Fortunato, M., Alet, F.,
1010 Ravuri, S., Ewalds, T., Eaton-Rosen, Z., Hu, W., Merose, A., Hoyer, S., Holland, G.,
1011 Vinyals, O., Stott, J., Pritzel, A., Mohamed, S., & Battaglia, P. (2023). Learning
1012 skillful medium-range global weather forecasting. *Science*, 382(6677), 1416–
1013 1421. <https://doi.org/10.1126/science.adi2336>

1014 Lang, S., Alexe, M., Chantry, M., Dramsch, J., Pinault, F., Raoult, B., Clare, M. C. A.,
1015 Lessig, C., Maier-Gerber, M., Magnusson, L., Bouallègue, Z. B., Nemesio, A. P.,
1016 Dueben, P. D., Brown, A., Pappenberger, F., & Rabier, F. (2024). *AIFS - ECMWF's*
1017 *data-driven forecasting system* (arXiv:2406.01465). arXiv.
1018 <http://arxiv.org/abs/2406.01465>

1019 Lees, T., Reece, S., Kratzert, F., Klotz, D., Gauch, M., De Bruijn, J., Kumar Sahu, R., Greve,
1020 P., Slater, L., & Dadson, S. J. (2022). Hydrological concept formation inside long
1021 short-term memory (LSTM) networks. *Hydrology and Earth System Sciences*,
1022 26(12), 3079–3101. <https://doi.org/10.5194/hess-26-3079-2022>

1023 Li, L., Carver, R., Lopez-Gomez, I., Sha, F., & Anderson, J. (2023). *SEEDS: Emulation of*
1024 *Weather Forecast Ensembles with Diffusion Models*.
1025 <https://doi.org/10.48550/ARXIV.2306.14066>

1026 Machac, D., Reichert, P., & Albert, C. (2016). Emulation of dynamic simulators with
1027 application to hydrology. *Journal of Computational Physics*, 313, 352–366.
1028 <https://doi.org/10.1016/j.jcp.2016.02.046>

1029 Meyer, D., Grimmond, S., Dueben, P., Hogan, R., & Van Reeuwijk, M. (2022). Machine
1030 Learning Emulation of Urban Land Surface Processes. *Journal of Advances in*
1031 *Modeling Earth Systems*, 14(3), e2021MS002744.
1032 <https://doi.org/10.1029/2021MS002744>

1033 Mironov, D., & Helmert, J. (n.d.). *Parameterization of Lakes in NWP and Climate Models*.

1034 Muñoz-Sabater, J., Dutra, E., Agustí-Panareda, A., Albergel, C., Arduini, G., Balsamo, G.,
1035 Boussetta, S., Choulga, M., Harrigan, S., Hersbach, H., Martens, B., Miralles, D.
1036 G., Piles, M., Rodríguez-Fernández, N. J., Zsoter, E., Buontempo, C., & Thépaut,
1037 J.-N. (2021). ERA5-Land: A state-of-the-art global reanalysis dataset for land

1038 applications. *Earth System Science Data*, 13(9), 4349–4383.
1039 <https://doi.org/10.5194/essd-13-4349-2021>

1040 Nath, S., Lejeune, Q., Beusch, L., Seneviratne, S. I., & Schleussner, C.-F. (2022).
1041 MESMER-M: An Earth system model emulator for spatially resolved monthly
1042 temperature. *Earth System Dynamics*, 13(2), 851–877.
1043 <https://doi.org/10.5194/esd-13-851-2022>

1044 Nearing, G., Cohen, D., Dube, V., Gauch, M., Gilon, O., Harrigan, S., Hassidim, A., Klotz,
1045 D., Kratzert, F., Metzger, A., Nevo, S., Pappenberger, F., Prudhomme, C., Shalev,
1046 G., Shenzis, S., Tekalign, T. Y., Weitzner, D., & Matias, Y. (2024). Global prediction
1047 of extreme floods in ungauged watersheds. *Nature*, 627(8004), 559–563.
1048 <https://doi.org/10.1038/s41586-024-07145-1>

1049 Ozaki, Y., Tanigaki, Y., Watanabe, S., Nomura, M., & Onishi, M. (2022). Multiobjective
1050 Tree-Structured Parzen Estimator. *Journal of Artificial Intelligence Research*, 73,
1051 1209–1250. <https://doi.org/10.1613/jair.1.13188>

1052 Pappenberger, F., Ramos, M. H., Cloke, H. L., Wetterhall, F., Alfieri, L., Bogner, K.,
1053 Mueller, A., & Salamon, P. (2015). How do I know if my forecasts are better?
1054 Using benchmarks in hydrological ensemble prediction. *Journal of Hydrology*,
1055 522, 697–713. <https://doi.org/10.1016/j.jhydrol.2015.01.024>

1056 Popov, S., Morozov, S., & Babenko, A. (2019). *Neural Oblivious Decision Ensembles for*
1057 *Deep Learning on Tabular Data* (Version 2). arXiv.
1058 <https://doi.org/10.48550/ARXIV.1909.06312>

1059 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., &
1060 Prabhat. (2019). Deep learning and process understanding for data-driven Earth

1061 system science. *Nature*, 566(7743), 195–204. <https://doi.org/10.1038/s41586->
1062 019-0912-1

1063 Sener, O., & Koltun, V. (2018). *Multi-Task Learning as Multi-Objective Optimization*
1064 (Version 2). arXiv. <https://doi.org/10.48550/ARXIV.1810.04650>

1065 Shukla, J. (1981). Dynamical Predictability of Monthly Means. *Journal of the*
1066 *Atmospheric Sciences*, 38(12), 2547–2572. <https://doi.org/10.1175/1520->
1067 0469(1981)038<2547:DPOMM>2.0.CO;2

1068 Shwartz-Ziv, R., & Armon, A. (2021). *Tabular Data: Deep Learning is Not All You Need*.
1069 <https://doi.org/10.48550/ARXIV.2106.03253>

1070 Slater, L. J., Arnal, L., Boucher, M.-A., Chang, A. Y.-Y., Moulds, S., Murphy, C., Nearing,
1071 G., Shalev, G., Shen, C., Speight, L., Villarini, G., Wilby, R. L., Wood, A., & Zappa,
1072 M. (2023). Hybrid forecasting: Blending climate predictions with AI models.
1073 *Hydrology and Earth System Sciences*, 27(9), 1865–1889.
1074 <https://doi.org/10.5194/hess-27-1865-2023>

1075 Thorpe, A., Bauer, P., Magnusson, L., & Richardson, D. (2013). *An evaluation of recent*
1076 *performance of ECMWF's forecasts*. <https://doi.org/10.21957/HI1EEKTR>

1077 Van Katwyk, P., Fox-Kemper, B., Seroussi, H., Nowicki, S., & Bergen, K. J. (2023). A
1078 Variational LSTM Emulator of Sea Level Contribution From the Antarctic Ice
1079 Sheet. *Journal of Advances in Modeling Earth Systems*, 15(12), e2023MS003899.
1080 <https://doi.org/10.1029/2023MS003899>

1081 Viterbo, P. (2002). *Land_surface_processes*.

1082 Wesselkamp, M., Moser, N., Kalweit, M., Boedecker, J., & Dormann, C. F. (2022).
1083 *Process-guidance improves predictive performance of neural networks for*
1084 *carbon turnover in ecosystems*. <https://doi.org/10.48550/ARXIV.2209.14229>

1085 Zwart, J. A., Oliver, S. K., Watkins, W. D., Sadler, J. M., Appling, A. P., Corson-Dosch, H.
1086 R., Jia, X., Kumar, V., & Read, J. S. (2023). Near-term forecasts of stream
1087 temperature using deep learning and data assimilation in support of
1088 management decisions. *JAWRA Journal of the American Water Resources*
1089 *Association*, 59(2), 317–337. <https://doi.org/10.1111/1752-1688.13093>
1090

Seite 21: [1] hat formatiert	Marieke Wesselkamp	19.10.24 14:05:00
------------------------------	--------------------	-------------------

Schriftart: (Standard) Times New Roman, 12 Pt.

Seite 21: [1] hat formatiert	Marieke Wesselkamp	19.10.24 14:05:00
------------------------------	--------------------	-------------------

Schriftart: (Standard) Times New Roman, 12 Pt.

Seite 21: [1] hat formatiert	Marieke Wesselkamp	19.10.24 14:05:00
------------------------------	--------------------	-------------------

Schriftart: (Standard) Times New Roman, 12 Pt.

Seite 21: [1] hat formatiert	Marieke Wesselkamp	19.10.24 14:05:00
------------------------------	--------------------	-------------------

Schriftart: (Standard) Times New Roman, 12 Pt.

Seite 21: [1] hat formatiert	Marieke Wesselkamp	19.10.24 14:05:00
------------------------------	--------------------	-------------------

Schriftart: (Standard) Times New Roman, 12 Pt.

Seite 21: [1] hat formatiert	Marieke Wesselkamp	19.10.24 14:05:00
------------------------------	--------------------	-------------------

Schriftart: (Standard) Times New Roman, 12 Pt.

Seite 21: [1] hat formatiert	Marieke Wesselkamp	19.10.24 14:05:00
------------------------------	--------------------	-------------------

Schriftart: (Standard) Times New Roman, 12 Pt.

Seite 21: [2] hat gelöscht	Marieke Wesselkamp	19.10.24 12:46:00
----------------------------	--------------------	-------------------

x.....

Seite 21: [2] hat gelöscht	Marieke Wesselkamp	19.10.24 12:46:00
----------------------------	--------------------	-------------------

x.....

Seite 21: [2] hat gelöscht	Marieke Wesselkamp	19.10.24 12:46:00
----------------------------	--------------------	-------------------

x.....

Seite 21: [2] hat gelöscht	Marieke Wesselkamp	19.10.24 12:46:00
----------------------------	--------------------	-------------------

x.....