



# 1 **Machine Learning for numerical weather and climate modelling:** 2 **a review**

3 Catherine O. de Burgh-Day<sup>1</sup> & Tennessee Leeuwenburg<sup>1</sup>

4 <sup>1</sup>The Bureau of Meteorology, 700 Collins St Docklands, Victoria, Australia

5 *Correspondence to:* Catherine O. de Burgh-Day ([catherine.deburgh-day@bom.gov.au](mailto:catherine.deburgh-day@bom.gov.au))

## 6 **Abstract.**

7 Machine learning (ML) is increasing in popularity in the field of weather and climate modelling. Applications range  
8 from improved solvers and preconditioners, to parametrisation scheme emulation and replacement, and recently even  
9 to full ML-based weather and climate prediction models. While ML has been used in this space for more than 25  
10 years, it is only in the last 10 or so years that progress has accelerated to the point that ML applications are becoming  
11 competitive with numerical knowledge-based alternatives. In this review, we provide a roughly chronological  
12 summary of the application of ML to aspects of weather and climate modelling from early publications through to the  
13 latest progress at the time of writing. We also provide an overview of key ML concepts and terms. Our aim is to  
14 provide a primer for researchers and model developers to rapidly familiarize and update themselves with the world of  
15 ML in the context of weather and climate models.

## 16 **1. Introduction**

17 Current state-of-the-art weather and climate models use numerical methods to solve equations representing the  
18 dynamics of the atmosphere and ocean on meshed grids. The grid-scale effects of processes that are too small to be  
19 resolved are either represented by parametrization schemes or are prescribed. These numerical weather and climate  
20 forecasts are computationally costly and are not amenable to transfer to specialized compute resources such as GPUs.  
21 One of the main approaches to improving forecast accuracy is to increase model resolution (reduced timestep between  
22 model increments and/or decreased grid spacing), but due to the high computational cost of this approach,  
23 improvements in model skill are hampered by the finite supercomputer capacity available. An additional pathway to  
24 improve skill is to improve the representation of sub grid-scale processes, however this is again a computationally  
25 costly exercise.

26 Machine learning is an increasingly powerful and popular tool. It has proven to be computationally efficient, as well  
27 as being an accurate way to model sub-grid scale processes. The term “Machine learning” (ML) was first coined by  
28 Arthur Samuel in 1952 to refer to a “field of study that gives computers the ability to learn without being explicitly  
29 programmed”<sup>1</sup>. Learning by example is the defining characteristic of ML.

---

<sup>1</sup> <http://infolab.stanford.edu/pub/voy/museum/samuel.html>, accessed 7<sup>th</sup> February 2023



30 The growing potential for ML in weather and climate modelling is being increasingly recognized by meteorological  
31 agencies and researchers around the world. The former is evidenced by the development of strategies and frameworks  
32 to better support the development of ML research, such as the Data Science Framework recently published by the Met  
33 Office in the UK<sup>2</sup>. The latter is made clear by the explosion in publications from academia, government agencies and  
34 private industry in this space, as demonstrated by the rest of this review.

35 There are established techniques and aspects of the weather and climate modeling lifecycle that would already be  
36 considered ML by many. For example, linear regression<sup>†3</sup>, principal component analysis, correlations, and the  
37 calculation of teleconnections can all be considered types of ML. Data Assimilation techniques could also be  
38 considered a form of ML. There are, however, other classes of ML (e.g. Neural Networks<sup>†</sup>, Decision Trees<sup>†</sup>, etc.)  
39 which are much less widely used within the weather and climate modelling space and have great potential to be of  
40 benefit. There is growing interest in, and increasingly effective application of, these ML techniques to take the place  
41 of more traditional approaches to modelling. The potential for ML in weather and climate modelling extends all the  
42 way from replacement of individual sub-components of the model (to improve accuracy and reduce computational  
43 cost) to full replacement of the entire numerical model.

44 While ML models are typically computationally costly during training, they can provide very fast predictions at  
45 inference<sup>†</sup> time, especially on GPU hardware. They often also avoid the need to have full understanding of the  
46 processes being represented and can learn and infer complex relationships without any need for them to be explicitly  
47 encoded. These properties make ML an attractive alternative to traditional parametrization, numerical solver, and  
48 modelling methods.

49 Neural Networks (NNs, explained further in Section 2.1) in particular are an increasingly favored alternative approach  
50 for representing sub-grid scale processes or replacing numerical models entirely. They consist of several  
51 interconnected layers of nonlinear nodes<sup>†</sup>, with the number of intermediate layers depending on the complexity of the  
52 system being represented. These nodes allow for the encoding of an arbitrary number of interrelationships between  
53 arbitrary parameters to represent the system, removing the need to explicitly encode these interrelationships into a  
54 parameterization or numerical model.

55 One challenge that must be overcome before there will be more widespread acceptance of ML as an alternative to  
56 traditional modelling methods is that ML is seen as lacking interpretability. Most ML models do not explicitly  
57 represent the physical processes they are simulating, although physics constrained ML is a new and growing field  
58 which goes some way to addressing this (see Section 6). Furthermore, the techniques available to gain insight into the  
59 relative importance and predictive mechanism of each predictor (i.e. the model outputs) are limited. In contrast,  
60 traditional models are usually driven by some understanding of the physical mechanisms and processes which are  
61 occurring. This makes it possible to more easily gain insight into what physical drivers could explain a given output.

---

<sup>2</sup> <https://www.metoffice.gov.uk/research/foundation/informatics-lab/met-office-data-science-framework>, accessed 7 February 2023

<sup>3</sup> Henceforth, the first occurrence of each term described in the glossary is marked with the symbol "†".



62 The “black box” nature of many ML approaches to parametrization makes them an unpopular choice for many  
63 researchers, and can be off-putting for decision makers, since, for example, explaining what went wrong in a model  
64 after a bad forecast can be more challenging if there are processes in the model which are not, and cannot, be  
65 understood through the lens of physics. However, increasing attention is being paid to the interpretability of ML  
66 models, and there are existing methods to provide greater insight into the way physical information is propagated  
67 through them (e.g., attention maps, which identify the regions in spatial input data that have the greatest impact on the  
68 output field, and ablation studies, which involve comparing reduced data sources and/or models to the original models  
69 that have full access to available data, to gain insight into the models).

70 As with their traditional counterparts, ML-based parametrizations and emulators are typically initially developed in  
71 single-column models, aquaplanet configurations, or otherwise simplified models. There are many examples of ML-  
72 based schemes which have been shown to perform well against benchmark alternatives in this setting, only to fail to  
73 do so in a realistic model setting. A common theme is that these ML schemes rapidly excite instabilities in the model  
74 because errors in the ML parametrization can push key parameters outside of the domain of the training data as the  
75 overall model is integrated forward in time, leading to rapidly escalating errors and to the model ‘blowing up’.  
76 Similarly, many ML-based full model replacements perform well for short lead times, only to exhibit model drift and  
77 a rapid loss of skill for longer lead times due to rapidly growing errors and the model drifting outside its training  
78 envelope.

79 In recent years, however, progress has been made in developing ML parametrizations which are stable within realistic  
80 models (i.e. not toy models, aquaplanets etc.), and ML-based full models which can run stably and skillfully to longer  
81 lead times. This is usually achieved through training the model on more comprehensive data, employing ML  
82 architectures which keep the model outputs within physically real limits, or imposing physical constraints or  
83 conservation rules within the ML architecture or training loss functions†.

84 There are still challenges and possible limitations to an ML approach to weather and climate modelling. In most cases,  
85 a robust ML model or parameterization scheme should be able to:

- 86 • remain stable in a full (i.e. non-idealized) model run,
- 87 • generalize to cases outside its training envelope,
- 88 • conserve energy and achieve the required closures.

89 Additionally, for an ML approach to be worthwhile it must provide one or more of the following benefits:

- 90 • For ML parametrization schemes:
  - 91 ○ a speedup of the representation of a sub-grid scale process vs. when run with a traditional
  - 92 parametrization scheme. This can make the difference between the scheme being cost-effective to
  - 93 run or not - when it is not cost-effective the process usually needs to be represented with a static
  - 94 forcing or boundary condition file,



- 95                   ○ a speedup of the model vs. when run with traditional parametrization schemes,
- 96                   ○ improved representation of sub-grid process(es) over traditional parameterization schemes, as
- 97                   measured by metrics appropriate to the situation,
- 98                   ○ improved overall accuracy/skill of the model when run with traditional parametrization schemes,
- 99                   ○ insight into physical processes not provided by current numerical models or theory.
- 100               • For full ML models:
- 101                   ○ a speedup of the model vs. an appropriate numerical model control,
- 102                   ○ improved overall accuracy/skill of the model vs. an appropriate numerical model control,
- 103                   ○ skillful prediction to greater lead times than an appropriate numerical model control,
- 104                   ○ insight into physical processes not provided by current numerical models or theory

105 Furthermore, in many cases of ML approaches to weather and climate modelling problems (particularly for full model  
106 replacement) the work is led by data scientists and ML researchers with limited expertise in weather and climate model  
107 evaluation. This can lead to flawed, misleading or incomplete evaluations. Hewamalage et al. (2022) have sought to  
108 rectify this problem by providing a guide to forecast evaluation for data scientists.

109 The scope of this review is deliberately limited to the application of ML within numerical weather and climate models  
110 or for their replacement. This is done to keep the length of this review manageable. ML has enormous utility for other  
111 aspects of the forecast value chain such as observation quality assurance, data assimilation, model output  
112 postprocessing, forecast/product generation, downscaling, impact prediction, decision support tools, etc. A review of  
113 the application of, and progress in, ML in these areas would be of great value but is outside the scope of this review  
114 and is left for future work. A brief introduction to key ML architectures and concepts, including suggested foundational  
115 reading, is also provided to aid readers who are unfamiliar with the subject.

116 The remainder of this review is structured as follows: In Section 2 a quick introduction to ML is provided. In Section  
117 3 ML use in sub-grid parametrization and emulation, along with tools and challenges specific to this domain, are  
118 covered. In Section 4 the application of ML for the partial differential equations governing fluid flow is reviewed. In  
119 Section 5 the use of ML for full model replacement or emulation is reviewed. In Section 6 the growing field of physics  
120 constrained ML models is introduced, and in Section 7 a number of topics tangential to the main focus of this review  
121 are briefly mentioned. In Section 8 a review of the history of, and progress in, ML outside of the fields of weather and  
122 climate science is presented. In Section 9 some practical considerations for the integration of ML innovations into  
123 operational and climate models are discussed, and Section 10 provides a summary. A Glossary of Terms is provided  
124 after the final Section to aid the reader in their understanding of key concepts and words.



## 125 2. A Quick Introduction to Machine Learning

126 While the scope of this paper is a review of ML work directly applicable to weather and climate modelling, an abridged  
127 introduction to some key fundamental ML concepts is provided here to aid the reader. Suggested starting points for  
128 interested readers include Chase et al (2022a, 2022b), Russell & Norvig (2021), Goodfellow et al. (2016), and Hastie  
129 et al. (2009).

130 This introductory section is a brief exposition of the concepts most central to this review. Definitions for this section  
131 can be found in the glossary.

132 The majority of ML methods which have found traction in weather and climate modelling were first developed in  
133 fields such as computer vision, natural language processing and statistical modelling. Few, if any, of the methods  
134 mentioned in this paper could be considered unique to weather and climate modelling, however, they have in many  
135 cases been modified to a greater or lesser extent to suit the characteristics of the problem. Furthermore, there is a trend  
136 towards increasingly customized architectures in this field as it matures.

137 In this review, the term algorithm refers to the mathematical underpinnings of a machine learning approach. By this  
138 definition, decision trees (DTs), NNs, linear regression and Fourier transforms are examples of algorithms. The two  
139 most relevant algorithms for this review are DTs and NNs.

140 The term architecture in machine learning refers to a specific way of utilizing an algorithm to achieve a modelling  
141 objective reliably. For example, the U-Net<sup>†</sup> architecture is a specific way of laying out a NN which has proven effective  
142 in many applications. The extreme gradient boosting decision tree<sup>†</sup> architecture is a specific way of utilizing DTs  
143 which has proven reliable and effective for an extraordinary number of problems and situations and is an excellent  
144 choice as a first tool to experiment with machine learning.

145 A major current focus of ML research is new architectures based on NNs. Research also continues into the algorithms  
146 themselves. Many other algorithms have been and continue to be employed in machine learning but are not central to  
147 this review.

148 A key point for ML researchers to be aware of is the critical importance of approaching model training carefully.  
149 There are many pitfalls which can result in underperformance, unexpected bias or misclassification. For instance,  
150 adversarial examples<sup>†</sup> can occur ‘naturally’, and systems which process data can be subject to adversarial attack<sup>†</sup>  
151 through the intentional supply of data designed to fool a trained network.

### 152 2.1. Introduction to Neural Networks

153 NNs can be regarded as universal function approximators (Hornik et al., 1989; see also Lu et al., 2019). Further, NN  
154 architectures can theoretically be themselves modelled as a very wide feed-forward<sup>†</sup> NN with a single hidden layer.  
155 A Fourier transform is another example of a function approximator, although it is not universal since not all functions  
156 are periodic. NNs can therefore be candidates for accurate modelling of physical processes.



157 ML models are typically introduced in the literature as being either classification<sup>†</sup> or regression<sup>†</sup> models, and either  
158 supervised<sup>†</sup> or unsupervised<sup>†</sup>.

159 The mathematical underpinning of a NN can be considered distinctly in terms of its evaluation<sup>†</sup> (i.e., output, or  
160 prediction) step and its training update step. The prediction step can be considered as the evaluation of a many-  
161 dimensional arbitrarily complex function.

162 The simplest NN is a single-input, single node network with a simple activation<sup>†</sup> function. A commonly used activation  
163 function is the sigmoid function, which helpfully compresses the range between 0 and 1. A classification model will  
164 employ a threshold to map the output into the target categories. A regression model seeks to optimize the output result  
165 against some target value for the function.

166 Complex NNs are built up from many individual nodes, which may have heterogenous activation functions and a  
167 complex connectome<sup>†</sup>. The forward pass<sup>†</sup>, by which inputs are fed into the network and evaluated against activation  
168 functions to produce the final prediction, uses computationally efficient processes to quickly produce the result.

169 The training step for a NN is far more complex. The earliest NNs were designed by hand rather than through  
170 automation. The training step applies a back-propagation<sup>†</sup> algorithm to apply adjustment factors to the weights<sup>†</sup> and  
171 biases<sup>†</sup> of each node based on the accuracy of the overall prediction from the network.

172 Training very large networks was initially impractical. Both hardware and architecture advances have changed this,  
173 resulting in the significant increase in application of NNs to practical problems. Most NN research explores how to  
174 utilize different architectures to train more effective networks. There is little research going into improving the  
175 prediction step as the effectiveness of a network is limited by its ability to learn rather than its ability to predict. Some  
176 research into computational efficiency is relevant to the predictive step. NNs can still be technically challenging to  
177 work with, and a lot of skill and knowledge are needed to approach new applications.

178 The major classes of NN architectures most likely to be encountered are:

- 179 • Token-sequence architectures, first applied to natural language processing, generation and translation;  
180 applicable to any time-series problem. Now also applied to image and video applications, and mixed-mode  
181 applications such as text-to-image or text-to-video
- 182 • Convolutional<sup>†</sup> architectures, first applied to image content recognition, which match the connectome of the  
183 network to the fine structure of images in hierarchical fashion to learn to recognize high-level objects in  
184 images
- 185 • Transformer architectures<sup>†</sup>, based on the attention mechanism<sup>†</sup> to provide a non-recurrent architecture which  
186 can be trained using parallelized training strategies. This allows larger models to be trained. Originally  
187 developed for sequence prediction and extended to image processed through vision transformer architectures.



188 **2.2. Introduction to Decision Trees**

189 DTs are a series of decision points, typically represented in binary fashion based on a simple threshold. A particular  
190 DT of a particular size maps the input conditions into a final 'leaf' node which represents the outcome of the decisions  
191 up to that point.

192 A random forest<sup>†</sup> (RF) is the composition of a large number of DTs assembled according to a prescribed generation  
193 scheme, which are used as an ensemble. A gradient boosted decision tree (GBDT) is built up sequentially, where each  
194 subsequent decision tree attempts to model the errors of the stack of trees built up thus far. This approach outperforms  
195 RFs in most cases.

196 The DT family of ML architectures are very easy to train and are very efficient. They are well documented in the  
197 public domain and in published literature. DTs are statistical in nature and are not capable of effectively generalizing  
198 to situations which are not similar to those seen during training. This can be an advantage when unbounded outputs  
199 would be problematic, however can lead to problems where an ability to produce out-of-training solutions is necessary.  
200 Additionally, current DT implementations require all nodes (of all trees in the case of RFs and GBDTs) to be held in  
201 memory at inference time, making them potentially memory heavy.

202 **3. Sub-grid parametrization and emulation**

203 Sub-grid scale processes in numerical weather and climate models are typically represented via a statistical  
204 parameterization of what the macroscopic impacts of the process would be on resolved processes and parameters.  
205 These are commonly referred to as parameterization schemes, and can be very complex and relatively computationally  
206 costly. For example, in the European Centre for Medium-Range Weather Forecast's (ECMWF) Integrated Forecasting  
207 System (IFS) model they account for about a third of the total computational cost of running the model (Chantry et al.  
208 2021b). They also require some understanding of the underlying unresolved physical processes. Examples of sub-grid  
209 scale processes which are typically currently parameterized in operational systems include gravity wave drag,  
210 convection, radiation, sub-grid scale turbulence, and cloud microphysics. As additional complexity (for example  
211 representation of aerosols, atmospheric chemistry, land surface processes, etc.) is added to numerical models, the  
212 computational cost will only increase.

213 ML presents an alternative approach to representing sub-grid scale processes, either by emulating the behavior of an  
214 existing parameterization scheme, emulating the behavior of sub-components of the scheme, by replacing the current  
215 scheme or sub-component entirely with an ML-based scheme, or by replacing the aggregate effects of multiple  
216 parameterization schemes with a single ML model.

217 ML emulation of existing schemes or sub-components has the advantage of maintaining the status quo within the  
218 model; no or minimal re-tuning of the model should be required since the ML emulation is trained to replicate the  
219 results of an already-tuned-for scheme. Because of this, the main benefit of this approach is that it reduces the  
220 computational cost of running the parameterization scheme. On the other hand, full replacement of an existing



221 parameterization scheme or sub-component with an ML alternative has the potential to be both computationally  
222 cheaper and also an improvement over the preceding scheme.

223 In the following subsections, a review of the literature on aspects of ML for the parametrization and emulation of sub  
224 grid-scale processes is presented.

### 225 **3.1. Early work on ML parametrization and ML emulations**

226 A popular target for applying ML in climate models is radiative transfer, since it is one of the more computationally  
227 costly components of the model. As such, many early examples of the use of ML in sub-grid parametrization schemes  
228 focus on aspects of this physical process. Chevallier et al. (1998) trained NNs to represent the radiative transfer budget  
229 from the top of the atmosphere to the land surface, with a focus on application in climate studies. They incorporated  
230 the information from both line-by-line and band models in their training to achieve competitive results against both  
231 benchmarks. Their NNs achieved accuracies comparable to or better than benchmark radiative transfer models of the  
232 time, while also being much faster computationally.

233 In contrast to the ML based scheme developed by Chevallier et al. (1998), which could be considered an entirely new  
234 parametrization scheme, Krasnopolsky et al. (2005) used NNs to develop an ML based emulation of the existing  
235 atmospheric longwave radiation parametrization scheme in the NCAR Community Atmospheric Model (CAM). The  
236 authors demonstrated speedups with the NN emulation of 50-80 times the original parameterization scheme.

237 Emulation of existing schemes has since then become a popular method for achieving significant model speedups. For  
238 example, Gettelman et al. (2021) investigated the differences between a GCM with the warm rain formation process  
239 replaced with a bin microphysical model (resulting in a 400% slowdown) and one with the standard bulk microphysics  
240 parameterization in place. They then replaced the bin microphysical model with a set of NNs designed to emulate the  
241 differences observed, and showed that this configuration was able to closely reproduce the effects of including the bin  
242 microphysical model, without any of the corresponding slowdown in the GCM.

### 243 **3.2. ML for coarse graining**

244 Coarse graining involves using higher resolution model or analysis data to map the relationship between smaller-scale  
245 processes and a coarser grid resolution. It can be used to develop parameterization schemes without explicitly  
246 representing the physics of smaller scale processes.

247 This has proven to be a popular method for developing ML-based parametrization schemes. Brenowitz & Bretherton  
248 (2018) used a near-global aqua planet simulation run at 4 km grid length to train a NN to represent the apparent sources  
249 of heat and moisture averaged onto 160 km<sup>2</sup> grid boxes. They then tested this scheme in a prognostic single column  
250 model and showed that it performed better than a traditional model in matching the behavior of the aqua planet  
251 simulation it was trained on. Brenowitz & Bretherton (2019) built on this work by training their NN on the same global  
252 aqua-planet 4 km simulation, but then embedded this scheme within a coarser resolution (160 km<sup>2</sup>) global aqua planet  
253 General Circulation Model (GCM). Embedding NNs within GCMs is challenging because feedbacks between NN and





254 GCM components can cause spatially extended simulations to become dynamically unstable within a few model days.  
255 This is due to the inherently chaotic nature of the atmosphere in the GCM responding to inputs from the NN which  
256 cause rapidly escalating dynamical instabilities and/or violate physical conservation laws. The authors overcame this  
257 by identifying and removing inputs into the NN which were contributing to feedbacks between the NN and GCM  
258 (Brenowitz et al. 2020), and by including multiple time steps in the NN training cost function. This resulted in stable  
259 simulations which predicted the future state more accurately than the coarse resolution GCM without any  
260 parametrization of subgrid-scale variability, however the authors do observe that the mean state of their NN-coupled  
261 GCM would drift, making it unsuitable for prognostic climate simulations.

262 Rasp et al. (2018) trained a Deep Neural Network<sup>†</sup> (DNN) to represent all atmospheric subgrid processes in an  
263 aquaplanet climate model by learning from a multiscale model in which convection was treated explicitly. They then  
264 replaced all sub-grid parameterizations in an aquaplanet GCM with the DNN, and allowed it to freely interact with  
265 the resolved dynamics and the surface-flux scheme. They showed that the resulting system was stable and able to  
266 closely reproduce not only the mean climate of the cloud-resolving simulation but also key aspects of variability in  
267 prognostic multiyear simulations. The authors noted that their decision to use DNNs was a deliberate one, because  
268 they proved more stable in their prognostic simulations than NNs, and they also observed that larger networks achieved  
269 lower training losses. However, while Rasp et al. (2018) were able to engineer a stable model that produced results  
270 close to the reference GCM, small changes in the training dataset or input and output vectors quickly led to the NN  
271 producing increasingly unrealistic outputs and causing model blow-ups (Rasp 2020). Consistent with this, Brenowitz  
272 & Bretherton (2019) report that they were unable to achieve the same improvements in stability with increasing  
273 network layers found by Rasp et al. (2018).

### 274 3.3. Overcoming instability in ML emulations and parametrizations

275 O’Gorman & Dwyer (2018) tackled the instabilities observed in NN- and DNN-based approaches to subgrid-scale  
276 parameterization by employing an alternative ML method; Random Forests (RFs; Breiman 2001; Tibshirani &  
277 Friedman 2001). The authors trained a RF to emulate the outputs of a conventional moist convection parameterization  
278 scheme. They then replaced the conventional parameterization scheme with this emulation within a global climate  
279 model, and showed that it ran stably and was able to accurately produce climate statistics such as precipitation  
280 extremes without needing to be specially trained on extreme scenarios. RFs consist of an ensemble of DTs, with the  
281 predictions of the RF being the average of the predictions of the DTs which in turn exist within the domain of the  
282 training data. RFs thus have the property that their predictions cannot go outside of the domain for their training data,  
283 which in the case of O’Gorman & Dwyer (2018) ensured conservation of energy and nonnegativity of surface  
284 precipitation (both critically important features of the moist convection parameterization scheme) were automatically  
285 achieved. A disadvantage of this method however is that it requires considerable memory when the climate model is  
286 being run to store the tree structures and predicted values which make up the RF.

287 Yuval & O’Gorman (2020) extended on the ideas in O’Gorman & Dwyer (2018), switching from emulation of a single  
288 parametrization scheme to emulation of all atmospheric sub grid processes. They trained an RF on a high-resolution



289 three-dimensional model of a quasi-global atmosphere to produce outputs for a course-grained version of the model,  
290 and showed that at course resolution the RF can be used to reproduce the climate of the high-resolution simulation,  
291 running stably for 1000 days.

292 There are some drawbacks to a RF approach compared to a NN approach however; namely that NNs provide the  
293 possibility for greater accuracy than RFs, and also require substantially less memory when implemented. Given that  
294 GCMs are already memory intensive this can be a limiting factor in the practical application of ML parametrization  
295 schemes. Furthermore, there is the potential to implement reduced precision NNs on Graphics Processing Units  
296 (GPUs) and Central Processing Units (CPUs) which still achieve sufficient accuracy, leading to substantial gains in  
297 computational efficiency. Motivated by these considerations, Yuval et al. (2021) trained a NN in a similar manner to  
298 how the RF in Yuval & O’Gorman (2020) was trained, using a high resolution aqua-planet model and aiming to coarse  
299 grain the model parameters. They overcame the model instabilities observed to occur in previous attempts to use NNs  
300 for this process by wherever possible training to predict fluxes and sources and sinks (as opposed to the net tendencies  
301 predicted by the RF in Yuval & O’Gorman (2020)), thus incorporating physical constraints into the NN  
302 parametrization. The authors also investigated the impact of reduced precision in the NN, and found that it had little  
303 impact on the simulated climate.

#### 304 **3.4. From aquaplanets to realistic land-ocean simulations**

305 All of the studies discussed in this section so far which were tested in a full GCM have used aqua planet simulations.  
306 Han et al. (2020) broke away from this trend by developing a Residual Neural Network (ResNet) based  
307 parametrization scheme which emulated the moist physics processes in a realistic land-ocean simulation. Their  
308 emulation reproduced the characteristics of the land-ocean simulation well, and was also stable when embedded in  
309 single column models.

310 Mooers et al. (2021) represents a subsequent example of an ML emulation of atmospheric fields with realistic  
311 geographical boundary conditions, where the authors developed feed-forward DNNs to super-parametrize sub grid-  
312 scale atmospheric parameters and forced a realistic land surface model with them. Super-parametrization is distinct  
313 from traditional parameterization in that it relies on solving (usually simplified) governing equations for sub grid-scale  
314 processes rather than heuristic approximations of these processes. They employed automated hyperparameter  
315 optimization<sup>†</sup> to investigate a range of neural network architectures across ~250 trials, and investigated the statistical  
316 characteristics of their emulations. While the authors found that their DNNs had a less good fit in the tropical marine  
317 boundary layer, attributable to the DNN struggling to emulate fast stochastic signals in convection, they also reported  
318 good skill for signals on diurnal to synoptic timescales.

319 Brenowitz et al. (2022) sought to address the challenge of emulating fast processes. They used FV3GFS (Zhou et al.,  
320 2019; Harris et al., 2021; a compressible atmospheric model used for operational weather forecasts by the US National  
321 Weather Service) with a simple cloud microphysics scheme included to generate training data and used this to train a  
322 selection of ML models to emulate cloud microphysics processes, including fast phase changes. They emulated



323 different aspects of the microphysics with separate ML models chosen to be suitable to each task. For example, simple  
324 parameters were trained with single-layer NNs, while parameters which are more complex spatially were trained with  
325 RNNs (e.g., rain falls downwards and not upwards, so it is sequential in timesteps through the atmosphere – a feature  
326 which can be represented by an RNN). They then embedded their ML emulation in FV3GFS. They found that their  
327 combined ML simulation performed skillfully according to their chosen metrics, but had excessive cloud over the  
328 Antarctic Plateau.

329 All of these studies, however, did not test their parameterizations in prognostic long-term simulations.

### 330 **3.5. Testing with prognostic long-term simulations**

331 A barrier to achieving stable runs with minimal model drift with ML components is the fact that generic ML models  
332 are not designed to conserve quantities which are required to be conserved by the physics of the atmosphere and ocean.  
333 Beucler et al. (2019) proposed and tested two methods for imposing such constraints in a NN model: (1) constraining  
334 the loss function or (2) constraining the architecture of the network itself. They found that their control NN with no  
335 physical constraints imposed performed well, but did so by breaking conservation laws, bringing into question the  
336 trustworthiness of such a model in a prognostic setting. Their constrained networks did however generalize better to  
337 unforeseen conditions, implying they might perform better under a changing climate than unconstrained models.

338 Chantry et al. (2021b) trained a NN to emulate the non-orographic gravity wave drag parameterization in the ECMWF  
339 IFS model (specifically cycle 45R1, ECMWF, 2018) and were able to run stable, accurate simulations out to 1 year  
340 with this emulation coupled to the IFS. While the authors note that RFs have been shown to be more stable (e.g.,  
341 O’Gorman & Dwyer (2018) and Yuval & O’Gorman (2020), as described above, and Brenowitz et al. (2020)), they  
342 chose to focus on NNs since they have lower memory requirements and therefore promise better theoretical  
343 performance. The authors assessed the performance of their emulation in a realistic GCM by coupling the NN with  
344 the IFS, replacing the existing non-orographic gravity wave drag scheme, and performed 120 hour, 10 day, and 1 year  
345 forecasts at ~25 km resolution in a variety of model configurations. The authors showed that their emulation was able  
346 to run stably when coupled to the IFS for seasonal timescales, including being able to reproduce the descent of the  
347 Quasi-biennial Oscillation (QBO). Interestingly, while the authors initially aimed to ensure momentum conservation  
348 in a manner similar to Beucler et al. (2021), they found that this constraint led to model instabilities and that a better  
349 result was achieved without it. One possible explanation for this is that Beucler et al. (2021) assessed their NNs in an  
350 aquaplanet setting. Nonetheless, Chantry et al. (2021b) noted that since their method was not identical to Beucler et  
351 al. (2021), improved stability could potentially be achieved by following their method more precisely. The  
352 computational cost of the NN emulation developed by Chantry et al. (2021b) was found to be similar that of the  
353 existing parametrization scheme when run on CPUs, but was faster by a factor of 10 when run on GPUs due to the  
354 reduction in data transmission bottlenecks.

355 The first study to successfully run stable long-term climate simulations with ML parametrizations was Wang et al.  
356 (2022a), who extended on the work of Han et al. (2020) by constructing a Residual DNN<sup>†</sup> (ResDNN) to emulate moist



357 physics processes. They used the residual connections from Han et al. (2020) to construct DNNs with good nonlinear  
358 fitting ability, and filtered out unstable NN parametrizations using a trial-and-error analysis, resulting in the best  
359 ResDNN set in terms of accuracy and long-term stability. They implemented this scheme in a GCM with realistic  
360 geographical boundary conditions and were able to maintain stable simulations for over 10 years in an Atmospheric  
361 Model Intercomparison Project (AMIP)-style configuration. This was more akin to a hybrid ML-physics based model  
362 than a traditional GCM with ML-based parametrization, because rather than embedding the ResDNN in the model  
363 code, the authors used a NN-GCM coupling platform through which the NNs and GCMs could interact through data  
364 transmission. This is in contrast to the approach employed in the Physical-model Integration with Machine Learning<sup>4</sup>  
365 (PIML) project and Infero<sup>5</sup>, which are both described in Section 3.11. One advantage to this approach noted by the  
366 authors is that it allows for a high degree of flexibility in the application of the ML component, however is likely to  
367 be less efficient than a fully-embedded ML model, due to the potential for data transmission bottlenecks.

### 368 **3.6. Training with observational data**

369 An alternative to using more complex and/or higher resolution models for training data is to train using direct  
370 observational data. For example, Ukkonen & Mäkelä (2019) used reanalysis data from ERA5 and lightning  
371 observation data to train a variety of different types of ML models to predict thunderstorm occurrence; this was then  
372 used as a proxy to trigger deep convection. ML models assessed were logistic regression, RFs, GBDTs, and DNNs,  
373 with the final two showing a significant increase in skill over convective available potential energy (CAPE; a standard  
374 measure of potential convective instability). One of the challenges of accurately reproducing the large-scale effects of  
375 convection is correctly identifying when deep convection should occur within a grid cell. The authors proposed that  
376 an ML model such as those they assessed could be used as the “trigger function” which activates the deep convection  
377 scheme within a GCM.

### 378 **3.7. ML for super parameterization**

379 Revisiting the topic of super parametrized sub grid scale processes introduced above, the use of ML for this approach  
380 was investigated in depth by Chattopadhyay et al. (2020). The authors introduced a framework for NN-based super  
381 parameterization, and compared the performance of this method against NN-based traditional parameterization (i.e.,  
382 based on heuristic approximations of sub grid-scale processes) and direct super parameterization (i.e., explicitly  
383 solving for the sub grid-scale processes) in a chaotic Lorenz '96 (Lorenz 1996) system that had three sets of variables,  
384 each of a different scale. They found that their NN-based super parameterization outperformed direct super  
385 parameterization in terms of computational cost, and was more accurate than NN-based traditional parameterization.  
386 The NN-based super parameterization showed comparable accuracy to direct super parameterization in reproducing  
387 long-term climate statistics, but was not always comparable for short-term forecasting.

<sup>4</sup> <https://turbo-adventure-f9826cb3.pages.github.io/>, accessed 7<sup>th</sup> February 2023

<sup>5</sup> <https://infero.readthedocs.io/en/latest/>, accessed 7<sup>th</sup> February 2023



388 **3.8. Stochastic parametrization schemes**

389 A more recent approach to the representation of sub grid-scale processes is via stochastic parameterization schemes,  
390 which can represent uncertainty within the scheme. There has been less focus on replacing these schemes with ML  
391 alternatives than non-stochastic schemes, however some progress has been made. Krasnopolsky et al. (2013) used an  
392 ensemble of NNs to learn a stochastic convection parametrization from data from a high-resolution cloud resolving  
393 model. In this case, the stochastic nature of the parametrization was captured by the ensemble of NNs. Gagne et al  
394 (2020) took a different approach, investigating the utility of generative adversarial networks (GANs) for stochastic  
395 parametrization schemes in Lorenz '96 (Lorenz 1996) models. In this case, the GAN learned to emulate the noise of  
396 the scheme directly, rather than implicitly representing it with an ensemble. They described the effects of different  
397 methods to characterize input noise for the GAN, and the performance of the model at both weather and climate  
398 timescales. The authors found that the properties of the noise influenced the efficacy of training. Too much noise  
399 resulted impaired model convergence and too little noise resulted in instabilities within the trained networks.

400 **3.9. ML parametrization and emulation for ocean models**

401 ML approaches to parameterization of subgrid-scale processes is not limited to the atmosphere. Krasnopolsky et al.  
402 (2002) presented an early application of NN for the approximation of seawater density, the inversion of the seawater  
403 equation of state, and a NN approximation of the nonlinear wave-wave interaction. More recently, Bolton & Zanna  
404 (2019) investigated the utility of Convolutional Neural Networks (CNNs) for parametrizing unresolved turbulent  
405 ocean processes and subsurface flow fields. Zanna & Bolton (2020) then investigated both Relevance Vector  
406 Machines<sup>†</sup> (RVMs) and CNNs for parameterizing mesoscale ocean eddies. They demonstrate that because RVMs are  
407 interpretable, they can be used to reveal closed-form equations for eddy parameterizations with embedded  
408 conservation laws. The authors tested the RVM and CNN parameterizations in an idealized ocean model and found  
409 that both improved the statistics of the coarse resolution simulation. While the CNN was found to be more stable than  
410 the RVM, the advantage of the RVM is the greater interpretability of its outputs.

411 **3.10. ML for representing or correcting a sub-component of a parametrization scheme**

412 An alternative method to replacing or emulating an entire parametrization scheme or schemes with ML is to target the  
413 most costly or troublesome sub-components of the scheme, and either replace those or make corrections to them.

414 Ukkonen et al. (2020) trained NNs to replace gas optics computations in the RTE-RRTMGP (Radiative Transfer for  
415 Energetics and Rapid and accurate Radiative Transfer Model for General circulation models applications-Parallel;  
416 Pincus et al., 2019) scheme. The NNs were faster by a factor of 1-6, depending on the software and hardware platforms  
417 used. The accuracy of the scheme remained similar to that of the original scheme.

418 Meyer et al. (2022) trained a NN to account for the differences between 1D cloud effects in the European Centre for  
419 Medium Range Weather Forecasting (ECMWF) 1D radiation scheme ecRad and 3D cloud effects in the ECMWF  
420 SPARTACUS (SPeedy Algorithm for Radiative TrAnsfer through CloUd Sides) solver. The 1D cloud effects solver



421 within ecRad, Tripleclouds, is favored over the 3D SPARTACUS solver because it is five times less computationally  
422 expensive. The authors show that their NN can account for differences between the two schemes with typical errors  
423 between 20% and 30% of the 3D signal, resulting in an improvement in Tripleclouds' accuracy with an increase in  
424 runtime of approximately 1%. By accounting for the differences between SPARTACUS and Tripleclouds rather than  
425 emulating all of SPARTACUS, the authors were able to keep Tripleclouds unchanged within ecRad for cloud-free  
426 areas of the atmosphere, and utilize the NN 3D correction elsewhere.

### 427 **3.11. Bridging the gap between popular languages for ML and large numerical models**

428 A common toolset for researchers to develop and experiment with different ML approaches to problems is Python  
429 libraries such as pytorch, scikit-learn, tensorflow, keras, etc., or other dynamically-typed, non-compiled languages.  
430 In contrast, numerical weather models are almost universally written in statically-typed compiled languages,  
431 predominantly Fortran. To make use of ML emulations or parameterizations in the models thus requires that they be:

- 432 (1) treated as a separate model periodically coupled to the main model (as is done between atmosphere and ocean  
433 models for example), or
- 434 (2) be manually re-implemented in Fortran, or
- 435 (3) that the pre-existing libraries used are somehow be made accessible within the model code.

436 Wang et al. (2022a; mentioned already above) opted for method 1, developing what could be considered a hybrid ML-  
437 physics based model rather than a traditional GCM with ML-based parametrization. In their study, the authors used a  
438 NN-GCM coupling platform through which the NNs and GCMs could interact through data transmission. One  
439 advantage to this approach noted by the authors is that it allows for a high degree of flexibility in the application of  
440 the ML component, however, is likely to be less efficient than a fully-embedded ML model, due to the potential for  
441 data transmission bottlenecks. This framework was then formalized by Zhong et al. (2023).

442 There are many examples where method 2 was used, such as Rasp et al. (2018), Brenowitz & Bretherton (2018),  
443 Gagne et al. (2019) and Gagne et al. (2020). The obvious disadvantage of this approach is that every change to the  
444 ML model being used requires reimplementing in the Fortran, and if the aim is to test a suite of ML models, this  
445 approach becomes untenable. Furthermore, this approach poses greater technical barriers for scientists developing  
446 ML-based solutions for numerical model challenges, since they must be sufficiently proficient in Fortran to  
447 reimplement models in it, rather than using existing user-friendly Python toolkits.

448 A solution lying somewhere between methods 2 and 3 was developed by Ott et al. (2020), who developed a Fortran-  
449 Keras Bridge (FKB) library that facilitated the implementation of Keras-like<sup>†</sup> NN modules in Fortran, providing a  
450 more modular means to build NNs in Fortran code. This however did not fully overcome the drawbacks posed by  
451 method 2 on its own; implementation of layers in the Fortran is still necessary, and any innovations in the Python  
452 modules being used would need to be mirrored in the Fortran library.



453 Finally, method 3 is being tackled by the Met Office in the PIML<sup>6</sup> project, and by ECMWF with an application called  
454 Inero<sup>7</sup>. These projects both seek to develop a framework which can be used by researchers to develop ML solutions  
455 to modelling problems in Python, and then integrate them directly into the existing codebase of the physical model  
456 (e.g., the Unified model at the UK Met Office). The approach used is to directly expose the compiled code  
457 underpinning the Python modules within the physical model code.

#### 458 **4. Application of ML for the partial differential equations governing fluid flow**

459 The representation and solving of the partial differential equations (PDEs) governing the fluid flow and dynamical  
460 processes in the oceans and atmosphere can be considered the backbone of weather and climate models. The solvers  
461 used to find solutions to these equations are typically iterative, and must solve the dynamics-governing equations of  
462 their model on every timestep and at every grid point. There has been growing interest in using ML to facilitate  
463 speedups and computational cost reductions in the preconditioning and execution of these solvers. Preconditioners are  
464 used to reduce the number of iterations required for a solver to converge on a solution, and usually do so by inverting  
465 parts of the linear problem. Many earlier studies focused on using ML to select the best preconditioner and/or PDE  
466 solver from a set of possible choices (e.g. Holloway & Chen, 2007; Kuefler & Chen, 2008; George et al., 2008; Peairs  
467 & Chen, 2011; Huang et al., 2016; and Yamada et al., 2018). Ackmann et al. (2020) approached the preconditioner  
468 part of the system more directly, using a variety of ML methods to directly predict the pre-condition of a linear solver,  
469 rather than using a standard preconditioner. Rizzuti et al. (2019) focused on the solver, using ML to apply corrections  
470 to a traditional iterative solver for the Helmholtz equation. Going a step further, a number of studies have used ML to  
471 replace the linear solver entirely (Ladický et al., 2015; Yang et al., 2016; Tompson et al., 2017).

472 Representation of the fluid equations in a gridded model poses a challenge because of the inability to resolve fine  
473 features in their solution. This leads to the use of course-grained approximations to the actual equations, which aim to  
474 accurately represent longer-wavelength dynamics while properly accounting for unresolved smaller-scale features.  
475 Bar-Sinai et al. (2019) trained a NN to optimally discretize the PDEs based on actual solutions to the known underlying  
476 equations. They showed that their method is highly accurate, allowing them to integrate in time a collection of  
477 nonlinear equations in 1 spatial dimension at resolutions 4× to 8× coarser than was possible with standard finite-  
478 difference methods.

479 Building on this, Kochkov et al. (2021) developed a ML-based method to accurately calculate the time evolution of  
480 solutions to nonlinear PDEs which used grids an order of magnitude coarser than is traditionally required to achieve  
481 the same degree of accuracy. They used convolutional NNs to discover discretized versions of the equations (as in  
482 Bar-Sinai et al., 2019), and applied this method selectively to the components of traditional solvers most affected by  
483 coarse resolution, with each NN being equation specific. They utilized the property that the dynamics of the PDEs  
484 were localized, combined with the convolutional layers of their NN enforcing translation invariance<sup>†</sup>, to perform their

<sup>6</sup> <https://turbo-adventure-f9826cb3.pages.github.io/> accessed 7<sup>th</sup> February 2023

<sup>7</sup> <https://infero.readthedocs.io/en/latest/> accessed 7<sup>th</sup> February 2023



485 training simulations on small but high-resolution domains, making the training set affordable to produce. An  
486 interesting feature of their training approach, which is growing in popularity, was the inclusion of the numerical solver  
487 in the training loss function: the loss function was defined as the cumulative pointwise error between the predicted  
488 and ground truth values over the training period. In this way, the NN model could see its own outputs as inputs,  
489 ensuring an internally-consistent training process. This had the effect of improving the predictive performance of the  
490 model over longer timescales, in terms of both accuracy and stability. Finally, the authors demonstrated that their  
491 models produced generalizable properties (i.e., although the models were trained on small domains, they produced  
492 accurate simulations over larger domains with different forcing and Reynolds number). They showed that this  
493 generalization property arose from consistent physical constraints being enforced by their chosen method.

494 An alternative to using ML to discover discretized versions of the PDE equations is to instead use NNs to learn the  
495 evolution operator of the underlying unknown PDE, a method often referred to as a DeepONet<sup>†</sup>. The evolution operator  
496 maps the solution of a PDE forwards in time and completely characterizes the solution evolution of the underlying  
497 unknown PDE. Because it is operating on the PDE, it is scale invariant and so bypasses the restriction of other methods  
498 that must be trained for a specific discretization or grid scale. Interest in, and the degree of sophistication of,  
499 DeepONets has grown rapidly in recent years (e.g., Lu et al., 2019; Wu & Xiu, 2020; Bhattacharya et al., 2020; Li et  
500 al., 2020a; Li et al., 2020b; Li et al., 2020c; Nelsen & Stuart, 2021; Patel et al., 2021; Wang et al., 2021; Lanthaler et  
501 al. 2022), to the point where the method is showing promising speedups: 3x faster than traditional solvers in the case  
502 of Wang et al. (2021).

503 The application of ML to the solving of PDEs and the preconditioning of PDE solvers has been a fruitful avenue of  
504 research to date. It has led to innovations which have proven useful even outside of the immediate field (e.g., Pathak  
505 et al. 2022 adapted innovations from DeepONets to use in fully ML-based weather models - this is discussed further  
506 in the next Section). This is likely in part because there are many areas of engineering and science which are active in  
507 progressing relevant research, leading to a greater overall pace of innovation. ML-based PDE solvers and  
508 preconditioners have not yet been tested in a physical weather and climate model. There are few theoretical reasons  
509 this could not occur and, if effective, result in significant computational efficiencies for traditional physical model  
510 architectures. This poses an interesting avenue for further research.

## 511 **5. Numerical model replacement/emulation**

512 The shift from using ML to emulate or replace parametrization schemes to using ML to replace the entire GCM has  
513 been made plausible by the increasing volume of training data available. The focus in this section will be on the  
514 challenge of completely replacing a GCM with a ML model.

515 There has been a flurry of activity in the use of ML for nowcasting (e.g. Ravuri et al., 2021), however, since the focus  
516 of this review is on weather and climate applications, these studies will not be elaborated on.





517 **5.1. Early work – 1D deterministic models**

518 Work on the use of ML to predict chaotic time-domain systems initially focused on 1-D problems, including 1-D  
519 Lorenz systems (e.g. Karunasinghe & Liong, 2006; Vlachas et al., 2018). Of particular interest is Vlachas et al. (2018),  
520 who used Long Short-Term Memory Networks (LSTMs<sup>†</sup>), which are well-suited to complex time domain problems.  
521 The recent popularization of convolutional LSTMs, which can also incorporate spatial information, suggests that  
522 revisiting the application of LSTMs for the prediction of spatially resolved chaotic systems could prove fruitful.

523 **5.2. Moving to spatially extended deterministic ML-based models**

524 Replacing a GCM entirely with an ML alternative was first suggested and tested in a spatially-resolved global  
525 configuration by Dueben and Bauer (2018), although for this study they only sought to predict a single variable  
526 (geopotential height at 500 hPa) on a 6 degree grid. Scher (2018) trained a CNN to predict the next model state of a  
527 GCM based on the complete state of the model at the previous step (i.e., an emulator of the GCM). Since this work  
528 was intended to be a proof-of-concept, the authors used a highly simplified GCM with no seasonal or diurnal cycle,  
529 no ocean, no orography, a resolution of ~625 km in the horizontal, and 10 vertical levels. Nonetheless, their ML model  
530 showed impressive capabilities; it was able to predict the complete model state several timesteps ahead, and when run  
531 in an iterative way (i.e., by feeding the model outputs back as new inputs) was able to produce a stable climate run  
532 with the same climate statistics as the GCM, with no long-term drift (even though no conservation properties were  
533 explicitly built into the CNN). Scher & Messori (2019) then extended on this, but continued the proof-of-concept  
534 approach. They investigated the ability of NNs to make skillful forecasts iteratively a day at a time to a lead time of a  
535 few days for GCMs of varying complexity, and explored a combination of other factors, including number of training  
536 years, the effects of model retuning, and the impact of a seasonal cycle on NN model accuracy and stability.

537 Sønnderby et al. (2020) took a more targeted approach, developing a NN to produce probabilistic precipitation forecasts  
538 to a lead time of 8 hours on a 2 x 2 km resolution grid covering 7000 x 2500 km over the continental United States,  
539 with temporal resolution of 2 min and latency (execution time) in the order of seconds. The desired lead time is an  
540 input parameter and time-stepping is not used. The focus here was producing rapid high-resolution short-term forecasts  
541 of a single key variable. Weyn et al. (2019) also aimed to predict a limited number of variables but focused more on  
542 the NWP to medium range time domain. They trained a CNN to predict 500 hPa geopotential height and 300 to 700  
543 hPa geopotential thickness over the Northern Hemisphere to up to 14-days lead time, showing better skill out to 3  
544 days than persistence, climatology, and a dynamics-based barotropic vorticity model, but not better than an operational  
545 full-physics weather prediction model.

546 Weyn et al. (2020) then improved on this significantly, with a Deep U-Net style CNN trained to predict four variables  
547 (geopotential height at 500 and 1000 hPa, 300 to 700 hPa geopotential thickness, and 2 m temperature) globally to 14  
548 days lead time. A major innovation in this study was their use of a cubed-sphere grid, which minimized distortions  
549 for planar convolution algorithms while also providing closed boundary conditions for the edges of the cube faces.  
550 Additionally, they extended their previous work to include sequence prediction techniques, making skillful predictions



551 possible to longer lead times. Their improved model outperformed persistence and a coarse resolution comparator (a  
552 T42 spectral resolution version of the ECMWF IFS model, with 62 vertical levels and ~2.8 degree horizontal  
553 resolution) to the full 14 days lead time, but was not as skillful as a higher resolution comparator (a T63 spectral  
554 resolution version of the IFS model with 137 vertical levels and ~1.9 degree horizontal resolution) or the operational  
555 subseasonal-to-seasonal (S2S) version of the ECMWF IFS.

556 Inching slightly closer to being competitive with physical models, Rasp & Thuerey (2021) developed a ResNet DNN  
557 model trained to predict geopotential height, temperature and precipitation to 5 days lead time and assessed it against  
558 the same set of physical models as Weyn et al. (2020). Their model was close to as skillful as the T63 spectral  
559 resolution version of the IFS model, and had better skill to the 5 day lead time than Weyn et al. (2020).

560 Keisler (2022) took an ambitious step forward, training a Graph Neural Network<sup>†</sup> (GNN) model to predict 6 physical  
561 variables on 13 atmospheric levels on a 1-degree horizontal grid, which the authors claim is ~50-2000 times larger  
562 than the number of physical quantities predicted by the models in Rasp & Thuerey (2021) and Weyn et al. (2020).  
563 Their model worked by iteratively predicting the state of the 6 variables 6 hours into the future (i.e., the output of each  
564 model timestep was the input into the next timestep), to a total lead time of 6 days. The authors showed that their  
565 model outperformed both Rasp & Thuerey (2021) and Weyn et al. (2020) in the variables common to all three studies.  
566 They suggested that the gain in skill seen over previous studies was due to the use of more channels<sup>†</sup> of information,  
567 and the higher spatial and temporal resolution of their model. Finally, they showed that their model was more skillful  
568 than NOAA's GFS physical model to 6 days lead time, but not as skillful as ECMWF's IFS.

569 Lam et al. (2022) also used GNNs to build their ML-based weather and climate model, GraphCast. This model was  
570 the most skillful ML-based weather and climate model at the time of writing this review. While the first ML-based  
571 weather and climate model to claim to exceed the skill of a numerical model was Pangu-Weather (Bi et al., 2022;  
572 described in greater detail in the following subsection), GraphCast exceeded the skill of both the ECMWF  
573 deterministic operational forecasting system, HRES, and also Pangu-Weather. Furthermore, Lam et al. (2022) paid  
574 particular attention to evaluating their model and HRES against appropriate measures, and included existing model  
575 assessment scorecards from ECMWF to evaluate them. GraphCast capitalized on the ability of GNNs to model  
576 arbitrary sparse interactions by adopting a high-resolution multi-scale mesh representation of the input and output  
577 parameters. It was trained on the ECMWF ERA5 reanalysis archive to produce predictions of five surface variables  
578 and six atmospheric variables, each at 37 vertical pressure levels, on a 0.25° grid. It made predictions on a 6-hourly  
579 timestep and was run autoregressively to produce predictions to a 10-day lead time. The authors demonstrated that  
580 GraphCast was more accurate than HRES on 90.0% of the 2760 variable and lead time combinations they evaluated.

### 581 **5.3. Probabilistic ML-based models**

582 A common criticism of ML approaches to weather and climate prediction is the difficulty of representing uncertainty  
583 and extremes. For example, Watson (2022) argued that while there is now an abundance of examples of ML being  
584 used for model parameterization schemes, full model replacement, downscaling, and PDE solvers (much of which is



585 covered in this review), there are relatively few examples which address the question of how well ML approaches can  
586 reproduce extreme events and statistics. There are however a growing number of examples where this is now being  
587 considered, many of which fall into the category of full-model replacement.

588 Clare et al. (2021) tackled this challenge by training a NN to predict full probability density functions of geopotential  
589 height at 500 hPa and temperature at 850 hPa at 3 and 5 days lead time, producing a probabilistic forecast which was  
590 comparable in accuracy to Weyn et al. (2020).

591 Weyn et al. (2021) have also explored probabilistic ML predictions using an ensemble of DNNs similar to the one  
592 described in Weyn et al. (2020). The authors expanded the number of variables predicted from 4 to 6, and used initial  
593 condition perturbation methods and variations in atmospheric representation similar to those used in traditional  
594 ensemble prediction to generate the ensemble of DNN predictions. They generated 320-member ensembles (much  
595 larger than could be affordably achieved with a physics-based model) and produced forecasts to 6 weeks lead time -  
596 considerably longer than any comparable work to date. The skill of the ensemble mean of the system, a control  
597 member, and the full ensemble were assessed against the same metrics from the ECMWF sub-seasonal to seasonal  
598 (S2S) prediction system. Their ML ensemble model had lower skill than the S2S system at shorter lead times, but was  
599 comparable in skill at longer lead times.

600 Pathak et al. (2022) developed a weather model called FourCastNet, leveraging the work on DeepONets described in  
601 Section 4. In particular, the authors used a type of DeepONet called a Fourier Neural Operator (FNO). FourCastNet  
602 produced predictions of 20 variables (including challenging-to-predict variables such as surface winds and  
603 precipitation) on five vertical levels with 0.25 degree horizontal resolution, and had competitive skill against the  
604 ECMWF IFS to 1 week lead time. The high horizontal resolution of their model enabled it to resolve extreme events  
605 such as tropical cyclones and atmospheric rivers, and the speed of the model facilitated the generation of large  
606 ensembles (1,000's of members). This suggests that prediction of very extreme events may be possible. The authors  
607 make the ambitious claim that with additional resources and further development, they anticipate that FourCastNet  
608 could match the capabilities of current NWP models on all timescales and at all vertical levels of the atmosphere.

609 Bi et al. (2022) achieved a significant milestone with their model Pangu-Weather, the first ML-based model to perform  
610 better than the ECMWF IFS to a lead time of 7 days based on RMSE and Anomaly Correlation Coefficient (ACC)  
611 across several variables including geopotential height and temperature at 500 hPa. Pangu-Weather featured two major  
612 innovations over FourCastNet:

- 613 1. It used 3D (latitude, longitude and height) input grids trained against 3D output grids. This enabled different  
614 levels of the atmosphere to share information, which was not possible in FourCastNet, in spite of predicting  
615 variables on multiple atmospheric levels, because the levels were treated independently. In contrast, Pangu-  
616 weather adopted a 3D convolutional method the authors name the 3D Earth-specific transformer (3DEST),  
617 which enabled the flow of information both horizontally and vertically.
- 618 2. It was made up of a series of models trained with different prediction time gaps. The motivation for this was  
619 that, as noted by the authors, when the goal is to produce forecasts to 5 days (for example), but the timestep



620 of the basic forecast model is relatively short (e.g. 6 hours), many iterative executions of the model are  
621 required, with the errors of each iteration feeding onto the next. A shorter model timestep results in greater  
622 overall errors (due to more iterations being required to reach the final forecast lead time), and a longer model  
623 timestep reduces this error. Motivated by this, the authors trained several versions of their model to predict  
624 to different timesteps on a single iteration. The overall forecast to a given leadtime was then constructed  
625 using the longest possible timesteps. For example, for a 7-day forecast, a 24-hour forecast is iterated 7 times,  
626 whereas for a 23-hour forecast, a 6-hour forecast is iterated 3 times, followed by a 3-hour forecast 1 time,  
627 and 1-hour forecast 2 times. The authors noted that this strategy was not effective to multiweek or longer  
628 timescales; they reported that training the model with a 28-day timestep was difficult for example, and  
629 suggested that more powerful or complex ML methods would be required to achieve this.

630 As well as the relatively broad measures of RMSE and ACC, the authors assessed the ability of their system to  
631 represent the intensity and track of selected tropical cyclones, and explored the potential for producing useful ensemble  
632 forecasts. They found that Pangu-Weather predicted the tracks of the cyclones considered with a high degree of  
633 accuracy compared to the ECMWF IFS, however it underestimated cyclone intensity. The authors attributed this to  
634 the training data they used (ERA5) also underestimating cyclone intensity. To assess ensemble predictions, they  
635 perturbed the initial state of the system with Perlin noise vectors to produce a 100-member ensemble of forecasts and  
636 calculated the RMSE and ACC of the ensemble mean for selected variables. They note that the ensemble mean  
637 forecasts performs worse than a single deterministic forecast for shorter lead times (e.g., 1 day), but better for longer  
638 lead times. The authors did not, however, investigate the properties of the spread of the ensemble or its utility for  
639 probabilistic forecasts or predicting statistical extremes.

640 As already mentioned above, the skill of Pangu-Weather was exceeded by GraphCast, although Lam et al. (2022) only  
641 assessed GraphCast in deterministic setting. Nonetheless, there is nothing stopping GraphCast from being used in an  
642 ensemble mode, and the assessment of GraphCast presented by Lam et al. (2022) was much more comprehensive and  
643 exacting than the assessment of Pangu-Weather presented by Bi et al. (2022).

644 It should be noted that all of the major milestones and high-profile ML models described in this section have relied  
645 on reanalysis datasets produced by physics-based models. The provision of higher resolution and higher quality open  
646 datasets have the potential to drive progress in this area as much as, if not more than, improvements and further  
647 research into ML algorithms.

#### 648 **5.4. Moving to more extensible models**

649 As the effectiveness of ML approaches are increasingly demonstrated in the literature, additional factors become clear  
650 in considering these models for both research and application. In a research setting, the ability to readily perform  
651 transfer learning to new problems and reduce training costs will be significant in supporting adoption by other  
652 researchers.



653 This need for greater flexibility in both the input data sources and predictive outputs of ML weather and climate  
654 models was recognized by Nguyen et al. (2023), who developed a transformer architecture-based ML model called  
655 ClimaX. This model was designed as a foundational model, trained initially on datasets derived from the CMIP6  
656 (Eyring et al., 2016) dataset, and designed to be readily retrained to specific tasks using transfer learning. The authors  
657 demonstrated the skill of ClimaX against simpler ML models, and in some cases a numerical model (ECMWF IFS),  
658 for a variety of tasks including weather prediction, sub-seasonal prediction, climate scenario prediction, and climate  
659 downscaling. The authors showed that ClimaX was able to make skillful predictions in scenarios unseen during the  
660 initial CMIP6 training phase. Furthermore, ClimaX used novel encoding and aggregation blocks in its architecture to  
661 achieve much more affordable training compute costs than other ML weather and climate models such as Graphcast,  
662 Pangu-weather and FourCastNet.

#### 663 **5.5. Benchmark datasets for ML weather models**

664 Providing open benchmark data for machine learning challenges has been as transformational for the machine learning  
665 field as improved algorithms, the publication of papers, or improvements in hardware.

666 As the interest and activity in the use of ML as a potential alternative to knowledge-based numerical GCMs has grown,  
667 the need for consistent benchmarks for the intercomparison of ML-based models has become increasingly clear. Rasp  
668 et al. (2020) addressed this need with the introduction of WeatherBench. On this platform, the authors provided data  
669 derived from the ERA5 archive that has been simplified and streamlined for common ML use cases and use by a broad  
670 audience. They also proposed a set of evaluation metrics which facilitate direct comparison between different ML  
671 approaches, and provided baseline scores in these metrics for simple techniques such as linear regression, some deep  
672 learning models and some GCMs. Weyn et al. (2020) chose datasets and assessment metrics consistent with  
673 WeatherBench to facilitate intercomparison of results. Rasp & Thuerey (2021) directly used the benchmarks provided  
674 by WeatherBench in their assessment. They demonstrated that their model outperformed previous submissions to  
675 WeatherBench, highlighting its value as a tool to allow intercomparability of ML-based weather models. Other  
676 examples of studies using WeatherBench data and analysis methods are Clare et al. (2021) and Weyn et al. (2021).  
677 The parameters of a good benchmark dataset were further elucidated by Dueben et al. (2022), who provided an  
678 overview of the current status of benchmark datasets for ML in weather and climate in use in the research community  
679 and provided a set of guidelines for how researchers could build their own benchmark datasets.

680 At the time of writing this review, assessments of ML-based models had chiefly (but not exclusively) focused on  
681 simple statistics like globally-averaged RMSE, and not reported in detail on the degree to which they accurately  
682 captured specific processes such as cyclone formation, climate drivers such as the El Nino Southern Oscillation, or  
683 large scale structures such as the jetstreams. A useful contribution from the scientific community would be to better  
684 quantify and articulate a suite of tests and statistics that could form a 'report card' to provide better insight into the  
685 value of new ML models.



## 686 5.6. A hybrid approach

687 Arcomano et al. (2022) present an approach which straddles the theme of this section and that of the following section  
688 (physics-constrained ML models). Following Wikner et al. (2020), they used a numerical atmospheric GCM and a  
689 computationally-efficient ML method called reservoir computing in a hybrid configuration called Combined Hybrid-  
690 Parallel Prediction (CHyPP). Their hybrid model is more accurate than the GCM alone for most state variables to a  
691 lead time of 7-8 days. They also demonstrate the utility of their hybrid model for climate predictions with a 10-year  
692 long climate simulation, for which they showed that the hybrid model had smaller systematic errors and more realistic  
693 variability than the GCM alone.

## 694 5.7. ML for predicting ocean variables

695 More recently, greater attention has been paid to the application of ML to the ocean, particularly for seasonal to multi-  
696 year prediction. Initial work in this space focused on directly predicting key indices such as the NINO 3.4 index. For  
697 example, Ham et al. (2019) trained a CNN to produce skillful El Niño Southern Oscillation (ENSO) forecasts with a  
698 lead time of up to one and a half years. A limiting factor for the application of ML to ocean variables is the lack of  
699 availability of observational data for training. To overcome this, the authors used transfer learning<sup>†</sup> to train their model  
700 first on historical simulations, and then on a reanalysis from 1871 to 1973. Data from 1984 to 2017 was reserved for  
701 validation. Ham et al. (2021) improved on this by including information about the current season in the network inputs  
702 as one-hot vectors<sup>†</sup>. Including this seasonality information led to an overall increase in skill relative to the model in  
703 Ham et al. (2019), in particular for forecasts initiated in boreal spring, a season which is particularly difficult to predict  
704 beyond.

705 Kim et al. (2022) improved on the performance of the 2D CNNs used in Ham et al. (2019) and Ham et al. (2021) for  
706 predicting ENSO by instead using a convolutional LSTM network with a global receptive field<sup>†</sup>. The move to a larger  
707 (global) receptive field for the convolutional layers enabled the network to learn the large-scale drivers and precursors  
708 of ENSO variability, and the use of a recurrent<sup>†</sup> architecture (in this case LSTM) facilitated the encoding of long-term  
709 sequential features with visual attention<sup>†</sup>. This led to a 5.8% improvement of the correlation coefficient for Nino3.4  
710 index prediction and 13% improvement in corresponding temporal classification with a 12-month lead time compared  
711 to a 2D CNN.

712 Taylor & Feng (2022) moved from prediction of indices to spatial outputs, training a Unet-LSTM<sup>†</sup> model on ECMWF  
713 ERA5 monthly mean Sea Surface Temperature (SST) and 2-m air temperature data from 1950-2021 to predict global  
714 2D SSTs up to a 24-month lead time. The authors found that their model was skillful in predicting the 2019-2020 El  
715 Niño and the 2016-2017 and 2017-2018 La Niñas, but not for the 2015-2016 extreme El Niño. Since they did not  
716 include any subsurface information in their training data (in contrast to Ham et al. (2019) and Ham et al. (2021), who  
717 included ocean heat content), they concluded that subsurface information may have been relevant for the evolution of  
718 that event.



719 It is clear from the small number of (but rapidly evolving) studies in this space that there is great promise for the use  
720 of ML for seasonal and multi-year prediction of ocean variables, with many avenues to pursue to achieve potential  
721 skill gains.

## 722 **5.8. ML for climate prediction**

723 The literature on the use of ML for prediction on seasonal to climate timescales is still relatively sparse compared to  
724 its use for nowcasting and weather prediction. Some examples have been covered in previous sections, such as Weyn  
725 et al. (2021) on subseasonal to seasonal timescales in the atmosphere, and Ham et al. (2019), Ham et al. (2021), Kim  
726 et al. (2022) and Taylor & Feng (2022) on seasonal to multiyear timescales in the ocean. A major cause for this sparsity  
727 is that deep learning typically requires large training datasets, and the available observation period for the earth system  
728 is too short to provide appropriate training data for seasonal to climate timescales in most applications. In the  
729 subseasonal to seasonal end, this may be overcome by including more slowly-varying fields in the training (e.g. ocean  
730 variables), by designing models to learn the underlying dynamics which drive long-term variability, and by including  
731 more physical constraints on the models. On the climate end these same methods could be beneficial, as well as  
732 transfer learning, as is done in Ham et al. (2019), and data augmentation<sup>†</sup> techniques. Additionally, interest is  
733 increasing in the use of ML to predict weather regimes and large-scale circulation patterns, which may prove beneficial  
734 in informing seasonal and climate predictions (Nielsen et al., 2022).

735 With the growing maturity of the field of ML for weather and climate prediction, there is every reason to believe the  
736 challenges of prediction on seasonal to climate timescales can be overcome.

## 737 **6. Physics constrained ML models**

738 As has been briefly touched on in previous sections, a promising and increasingly popular method for improving the  
739 performance of ML applications in weather and climate modelling is to include physics-based constraints in the ML  
740 model design (e.g. Karpatne et al., 2017; de Bézenac et al., 2017; Beucler et al., 2019; Yuval et al., 2021; Beucler et  
741 al., 2021; Harder et al., 2022). This can be done through the overall design and formulation of the model, and through  
742 the use of custom loss functions which impose physically-motivated conservations and constraints.

743 An excellent review of the possible methods for incorporating physics constraints into ML models for weather and  
744 climate modelling, along with 10 case studies of noteworthy applications of these methods, is presented in Kashinath  
745 et al. (2021). The scope of Kashinath et al. (2021) is broad and includes studies not applied directly in the context of  
746 weather and climate modelling, but applicable to it. Rather than repeat the total of this summary here, the reader is  
747 directed to this review.

748 A class of physics-leveraged ML which has grown rapidly in popularity is Physics Informed Neural Networks  
749 (PINNs). These are discussed in Kashinath et al. (2021), but have also become a very active area of research since the  
750 publication of that review. A more up-to-date review of this class of NNs is presented by Cuomo et al. (2022), along  
751 with a review of other related Physics guided ML architectures.



752 While PINNs are an exciting and promising new NN architecture, they still face some challenges. For example, they  
753 have had little success simulating dynamical systems whose solution exhibits multi-scale, chaotic or turbulent  
754 behavior. Wang et al. (2022b) attributed this to the inability of PINNs to represent physical causality, and developed  
755 a solution by re-formulating the loss function of a PINN to explicitly account for physical causality during model  
756 training. They demonstrated that this modified PINN was able to successfully simulate chaotic systems such as a  
757 Lorenz system, and the Navier-Stokes equations in the turbulent regime; something which traditional PINNs were  
758 unable to do.

759 Nonetheless, recent work with PINNs has led to some interesting results for weather and climate simulation: Bihlo &  
760 Popovych (2022) used PINNs to solve the shallow-water equations on a rotating sphere, as a demonstration of their  
761 utility in a meteorological context, and Fuhg et al. (2022) developed a modified PINN to solve interval and fuzzy  
762 partial differential equations, enabling the solving of PDEs including uncertain parameter fields.

## 763 **7. Other applications of ML and considerations for the use of ML in Weather and Climate Models**

764 Aside from the most active areas of development in the use of ML in weather and climate models discussed in the  
765 sections above, there are a few areas of the literature worth mentioning that are adjacent to the main focus of this  
766 review. These topics are covered in the following subsections.

### 767 **7.1. Nudging**

768 Rather than replacing a component or components of a GCM with an ML alternative to gain skill improvements, Watt-  
769 Meyer et al. (2021) focused on using corrective nudging to reduce model biases and the errors they can introduce  
770 through feedbacks. The authors used RFs to learn bias-correcting tendencies from a hindcast nudged towards  
771 observations. They then coupled this RF to a prognostic simulation and attempted to correct the model drift with the  
772 learned nudging tendencies. While this simulation ran stably over the year-long test period and showed improvements  
773 in some variables, the errors in others were observed to increase. So far studies in this space seem to be limited to  
774 Watt-Meyer et al. (2021), however this method seems promising, so hopefully interest in developing this approach  
775 further will grow in the future.

### 776 **7.2. Object identification within models**

777 An alternative to achieving greater model accuracy through increasing resolution of the entire model grid is to develop  
778 techniques to identify critical systems and physical phenomena within the model, and embed higher resolution  
779 temporary subgrids within the larger GCM to more accurately simulate those processes. A key challenge to overcome  
780 to achieve this is automatically identifying key model features. For example, Mudigonda et al. (2017) investigated the  
781 feasibility of using a variety of NN architectures to identify storms, tropical cyclones and atmospheric rivers within  
782 model data, with promising results. A major limitation of this area of research is the frequent need for labelled datasets  
783 of the events being identified, which are currently quite limited. While there are approaches to this problem which





784 utilize unsupervised learning (i.e., learning without an objective function or labelled data), it is harder to achieve a  
785 meaningful result this way.

### 786 **7.3. Uncertainty quantification**

787 A common criticism of some ML models such as NNs is that it is difficult to represent the uncertainty of their outputs.  
788 Some examples of studies that have sought to overcome this have already been mentioned in Section 3.8, and there  
789 are other examples in the literature (e.g. Grigo & Koutsourelakis, 2019; Atkinson, 2020; Yeo et al., 2021; O'Leary et  
790 al., 2022), however it is nonetheless still a relatively underexplored aspect of ML models for physical systems. Psaros  
791 et al. (2022) suggest that this may be because they are also under-utilized within the broader deep learning community,  
792 and it is thus a developing field that is not universally trusted and understood yet. They also point out that the physical  
793 considerations inherent to ML applied to physical systems often make them more complicated and computationally  
794 expensive than standard ML applications, further disincentivizing the inclusion of uncertainty quantification in an  
795 already complex problem.

796 Only recently has attention to this aspect of ML become sufficient to motivate the collection of methods into a  
797 consistent framework, a good example of which is the aforementioned Psaros et al. (2022), who presented a  
798 comprehensive review of the methods for quantifying uncertainty in NNs and provided a framework for applying  
799 these methods.

800 A related topic which is facing similar challenges is the question of explainability of ML approaches; often there is  
801 value in understanding the relative roles and importance of predictors in an ML model, or the relative significance of  
802 different regions of the predictor data. Flora et al. (2022) provide a good overview of approaches to this and compare  
803 their relative drawbacks and benefits.

### 804 **7.4. GPUs and specialized compute resources**

805 GPUs and TPUs are specialized hardware which are well suited to highly parallelizable matrix operations, ideal for  
806 solving neural network operations. TPUs have been developed specifically for deep learning applications. Both GPUs  
807 and TPUs are likely to be available on many of the next generation of supercomputers, but much of the current Fortran-  
808 based numerical weather and climate model infrastructure cannot be run on them in their current state. Data  
809 bottlenecks also exist between the GPUs (which have their own on-board memory) and the main memory accessible  
810 to the CPU. While efforts are underway to make numerical and climate models better suited to GPUs, for example  
811 with the development of LFRic (Adams et al. 2019), the new weather and climate modelling system being developed  
812 by the UK Met Office to replace the existing Unified Model (Walters et al. 2017), there is still a long way to go before  
813 entire weather and climate models can be reliably run on GPU or other specialized compute architectures. At the same  
814 time, some neural network designs are aimed squarely at the partial differential equation solving at the core of  
815 numerical methods. Since neural network evaluation utilizes simpler mathematical operations than current PDE  
816 solvers, they offer the prospect of significant computational advantages on non-specialized (i.e., CPU) hardware.



## 817 **8. Key papers from computer science**

818 This section provides a brief perspective on weather and climate modelling from the computer science domain, and  
819 aims to provide the earth system scientist with a short list of the main relevant innovations in computer science.

820 As was noted in Section 1, ML models are often regarded as black-boxes, largely because of the design of many  
821 prominent ML systems. In principle, it is not quite right to refer to a model as "a machine learning model", in the sense  
822 that either a physical model or a NN could undergo a training cycle to determine optimal parameter values. The  
823 parameters of a neural network are its weights and biases, whereas the parameters of a physical model are physical  
824 variables and constants. The essence of ML is the level of automation involved. Even in typical ML models such as  
825 large NNs, the model architecture is typically specified manually by the data scientist or physical scientist involved.  
826 The automated derivation of model architecture and composition is not yet mature for large models, although it is  
827 explored through evolutionary programming techniques whereby the learning of architecture as well as  
828 parameterization is automated.

829 As such, the goal of a "ML weather/climate model" (either a full model or an augmented numerical model) will likely  
830 be achieved using multiple model types and architectures, in a complex fashion.

831 Some models from the computer science domain also have a more fundamental probabilistic or statistical underpinning  
832 than typical weather and climate models in that the model state variables comprise a probability distribution or degree  
833 of confidence. Whereas a typical weather or climate model derives its probability outputs from an ensemble of  
834 perturbed members, an alternative approach could be taken whereby each part of the belief state<sup>†</sup> of the model is a  
835 distribution or likelihood, built up either empirically or by fitting a gaussian or other known distribution.

836 As such, ML may be applied to statistical models, process-based models, Bayesian models or physical models.  
837 Nonetheless, current directions in ML are focused on very large or deep NNs which rely both on the universal  
838 approximation theorem and practical experimentation to capture a prediction function without needing to explicitly  
839 represent the processes being modeled. In a conceptually similar fashion to how a Fourier decomposition can represent  
840 any wavelike function, the universal approximation theorem establishes that a NN may approximate any function,  
841 subject to its size and the required degree of accuracy (Hornik, Stinchcome and White 1989). Deep learning has been  
842 highly effective in approaching many problems, but many limitations are acknowledged, as evidenced by the current  
843 widespread focus on trustworthy computing and efforts towards explainable ML systems.

### 844 **8.1. A Selected Timeline of Machine Learning:**

845 In this section a selected history of ML developments relevant to weather and climate modelling is presented. The  
846 purpose of this is to give the reader some perspective on the timelines of the innovations in this space, and provide  
847 some additional context for the current state of the field.

- 848 • 1943: First development of a mathematical model of a human neuron by McCulloch & Pitts (1943)
- 849 • 1950: The Turing Test (originally called The Imitation Game) was developed (Turing, 1950)



- 850       • 1956: Arthur Samuel developed the first automated checkers program (one of the earliest examples of using  
851 machine learning to exceed the human skill of the developer)
- 852       • 1974: Development of backpropagation as a sound way of updating a parameterization based on model errors  
853 (Werbos, 1974, 1990). This work showed how to apply machine learning to artificial neuron models (i.e.,  
854 NNs).
- 855       • 1989: Hornik et al. (1989) proved that NNs can effectively approximate any function. This had theoretical  
856 implications for speedups to equation solvers which could theoretically be solved with orders of magnitude  
857 greater efficiency
- 858       • 1996: The LSTM architecture was developed, addressing the vanishing gradient problem (Hochreiter &  
859 Schmidhuber, 1996). LSTMs and their successors would revolutionize language modelling, speech  
860 recognition, translation, and text-to-speech synthesis. They can be used in any tokenized sequence modelling  
861 and have many applications in time-series prediction.
- 862       • 1998: The release of the Modified National Institute of Standards and Technology (MNIST) database spurred  
863 a generation of machine learning research into novel methods, coinciding with and stimulating major research  
864 and application of new techniques including DTs and NN approaches (LeCun et al., 1998). This milestone  
865 demonstrates how the release of open data, in an accessible way which supports machine learning patterns,  
866 can accelerate progress in a research field. Many people's first experience with machine learning is to train a  
867 model for MNIST prediction.
- 868       • 2001: Breiman (2001) developed RFs. This CPU-efficient algorithm found widespread application, and  
869 evolutions of this approach remain relevant today. This approach was extended to gradient-boosting by  
870 Freidman (2001).
- 871       • 2009: ImageNet data set released (Deng et al., 2009). Called "the data that transformed AI research", this  
872 milestone shows again how a key dataset release designed to support machine learning research can  
873 accelerate an entire field. This was also an important moment in the AI industry better recognizing the  
874 importance of data quality as much as algorithm complexity to AI performance.
- 875       • 2014: Goodfellow et al., (2020) introduced GANs (note that while this paper was ultimately published in  
876 2020, a version of the paper was available on ArXiv in 2014<sup>8</sup>). GANs introduce the concept of a generator<sup>†</sup>  
877 model, and a discriminator model which aims to distinguish predictions from observations. While the  
878 generator learns the predictive function, its objective function is the discriminator. When the discriminator  
879 can no longer distinguish the prediction from an observation, the overall model has been optimized. This has  
880 applications in producing highly realistic gridded outputs.

---

<sup>8</sup> <https://doi.org/10.48550/arXiv.1406.2661>, accessed 7<sup>th</sup> February 2023



- 881       • 2015: Ronneberger et al. (2015) introduced U-Nets, which are useful for image segmentation. Image  
882       segmentation has application in locating weather features such as fronts, cyclones, convection etc.
- 883       • 2015: He et al. (2016) introduced CNNs with skip-connections to improve performance and allow attention-  
884       like mechanisms.
- 885       • 2017: Vaswani et al. (2017) provided a significant step forward for sequence transduction modelling (i.e.,  
886       language translation) using an 'attention' mechanism which allows models to consider contextual factors in  
887       sequence prediction. Attention mechanisms let a model draw from the state at any preceding point along the  
888       sequence, allowing long sequences to be learned. This approach is replacing recurrent architectures for many  
889       applications.
- 890       • 2020: Dosovitskiy et al. (2021) applied and extended attention-based techniques to image processing with  
891       reduced computational requirements.

892       This history shows the degree and rate of research into processing images, text and other sequences based on semantic  
893       understanding of content, but does not demonstrate capturing physical processes as a core element. Advances in the  
894       weather and climate modelling domain have a more explicit goal of properly portraying real physical processes.  
895       Bringing these concepts together promises to uplift capability in both fields.

## 896       **9. Practical Perspectives on Machine Learning for Weather and Climate Models**

897       A major driver of research into, and improvement of, weather and climate models is increasing the skill of operational  
898       forecast systems around the world, and increasing the accuracy and trustworthiness of climate projections. Therefore,  
899       an important consideration for ML in the context of weather and climate models is the need for it to ultimately be  
900       integrated into a complete predictive system with practical application for forecasting or climate projections.

901       The research findings covered in this review, however compelling, are yet to make major changes to operational  
902       modelling systems, or standard climate projections. There are a number of challenges (which are discussed further in  
903       Section 10) that need to be overcome to leverage and capitalize on the new capabilities which have been developed in  
904       research settings.

905       This section summarizes some practical considerations the research community may wish to be aware of.

906       A major function of operational meteorological services is to inform of future conditions, largely for managing risk  
907       or optimizing benefits. A conservative approach is taken, which is to say, the utmost premium is put on accuracy,  
908       resilience, reliability, and solid scientific foundation. There is often a large gap between a research finding and the  
909       ability to successfully integrate an innovative new approach into a major model upgrade or scientific configuration  
910       upgrade. Understanding when to invest effort in bringing a research innovation into a major model or scientific  
911       configuration upgrade is a significant challenge requiring a great deal of time and effort. This is also true of the  
912       transition into operations for operational agencies.



913 One pathway to adoption of weather and climate models that use ML could be the development of limited-scope  
914 models optimized for one or a few parameters. Early effectiveness of limited-purpose ML models provide the ability  
915 to augment existing services without disruption. A risk associated with this approach is that inconsistencies between  
916 predictions may arise from independent forecasts from a collection of limited-scope models, leading to confusion from  
917 users and an erosion of trust.

918 Such an approach could be considered to entirely replace current model systems at the cost of consistency and model  
919 fidelity. Operational development is typically more incremental, however. It is more likely that progress will be made  
920 in achievable increments, along the evolving technical frontier. However promising and fascinating as a research  
921 direction, full model replacement is currently not mature enough for an operational system.

922 What may currently be operationally feasible includes parameterization scheme replacement, solver replacement,  
923 super-resolution, new approaches to data assimilation of novel observation sources, and both pre- and post-processing  
924 applications (although of course not all these applications have been covered in this review).

925 It is expected that the research into, and application of, ML methods will represent an increasing proportion of weather  
926 and climate model research, with increasingly sophisticated and skillful model components finding their way into  
927 major model releases over the coming years. These components are appealing for both computational and model skill  
928 reasons, and are expected to be highly promising avenues of research.

## 929 **10. Conclusions**

930 In this review we have presented a comprehensive survey of the literature on the use of ML in weather and climate  
931 modelling.

932 We have found that the ML models being most often explored include RFs, and NNs and DNNs (including ResNets  
933 and Deep U-Nets). We have also identified some recent innovations which have proven to be highly effective in the  
934 weather and climate modelling space, including DeepONets and variants thereof, and PINNs.

935 This review has demonstrated that ML is being successfully applied to many aspects of weather and climate modelling.  
936 We have presented examples from the literature of its application in (1) the emulation and replacement of sub-grid  
937 scale parametrizations and super-parametrizations, (2) preconditioning and solving of resolved equations, (3) full  
938 model replacement, and (4) a selection of other adjacent areas.

939 Nonetheless, there are still many challenges to overcome, including:

- 940 • addressing the instabilities excited in physical models due to the inclusion of ML components;
- 941 • increasing the ease of technical integration (in particular, Fortran compatibility);
- 942 • memory and computational concerns;



- 943       • representing a sufficient number of physical parameters and increasing physical and temporal resolution in  
944       ML-based weather and climate model implementations (which currently feature reduced fields and levels  
945       compared to physics-based numerical models);
- 946       • moving from a focus on individual parts of the earth system (i.e., the atmosphere, the ocean, the land surface  
947       etc.) to tackling the challenges associated with coupled models (i.e., where models of individual components  
948       of the earth system are coupled together). Increasingly, operational weather and climate models are coupled  
949       land-atmosphere-ocean-sea-ice models in order to more accurately represent the relevant timescales and  
950       processes in the earth system, and ML modelling efforts need to reflect this;
- 951       • the need for more good quality training data; and
- 952       • the practical challenges of integrating ML components or models into an operational setting.

953       This list provides a set of focus areas for future research efforts.

954       If the current trend in skill gains in full ML weather and climate models continues, it is possible they will eventually  
955       be considered viable alternatives to traditional numerical models. However, a more likely scenario is that ML  
956       components will replace an increasing number of physics-based model components, with models the near-term future  
957       being hybrid ML-physical models.

958       Some possible avenues through which increases in ML-based weather and climate model skill might be achieved is  
959       by operating at higher resolutions, resolving more processes which are implicit in the training data, or by undertaking  
960       experiments on synthetic data to address the paucity of real-world data.

961       Another benefit of ML approaches to weather and climate modeling is the relative computational cheapness of ML  
962       alternatives to current physics-based modelling systems. This has the potential to open the door to experiments that  
963       would not be feasible otherwise. For example, experiments requiring a very large ensemble would be more feasible  
964       with a computationally cheap ML approach.

965       The literature reviewed here indicates that 'out of the box' ML approaches and architectures are not effective when  
966       used in a weather and climate modelling context. Rather, ML architectures must be adapted to satisfy conservation of  
967       energy, represent physically realistic predictions and processes, and maintain good model stability. At the same time,  
968       computational and memory tractability must be maintained.

969       Advances in the sophistication, complexity and efficiency of ML architectures are being heavily invested in for many  
970       use cases in other disciplines and in the private sector (e.g., condition-action post estimation, text to video generation,  
971       stable diffusion/text to image, chatbots, facial recognition, semantic image decomposition, etc.). In order to capture  
972       the full benefits of ML for the weather and climate modelling domain, academic and operational agencies will need  
973       to continue to support research in this space.



974 Interest and progress in the application of ML to weather and climate modelling has been present for close to 30 years,  
975 and has begun to accelerate rapidly in the last few years. There is good reason to believe that ML as a tool will have  
976 transformational (and potentially highly fascinating, innovative, and offer great opportunity for further application.

### 977 **Machine Learning Glossary of Terms**

978 This glossary includes terms which the reader will come across frequently in machine learning literature for the  
979 weather and climate, as well as in machine learning literature generally. Most of these terms are used in this paper  
980 while others support further reading.

981 **Activation function.** The function which is used to multiply input values, add the bias and produce an output value  
982 from an individual node. Examples include linear, sigmoid and tanh.

983 **Adversarial attack.** The deliberate use of malicious data input in a real-world setting intended to cause a  
984 misclassification, underperformance or unexpected behaviours. Examples include emails designed to avoid spam  
985 filters, or images that have been modified to avoid recognition.

986 **Adversarial example.** A specialised input which results in a misclassification or underperformance of a predictive  
987 model. An example of this concept is an image which has had subtle noise added to it resulting in a copy of that image  
988 which is visually indistinguishable from the original, but which nonetheless causes a misclassification. The term  
989 'adversarial' is used to refer to the way the example fools the model and is not necessarily intended to convey the  
990 sense of malicious intent, although the term is often applied in that fashion. Adversarial examples demonstrate that  
991 machine learning models may be more brittle than expected based on ordinary training data alone. To increase model  
992 robustness, adversarial examples may be generated and added to the training set. Data augmentation techniques such  
993 as flipping, warping and adding noise (any many other techniques) are also used to generate additional training data  
994 to increase robustness and performance.

995 **Attention mechanism.** A complex mechanism to allow sequence prediction models to increase the importance of key  
996 terms within that sequence which may be nonlocal and modified in meaning according to the other terms of the  
997 sequence.

998 **API.** Application Programming Interface. A set of programming functions, methods or protocols by which to build  
999 and integrate applications. APIs may be "web" APIs or imported from software packages in which case they are more  
1000 often referred to as libraries.

1001 **Autoencoder.** A neural network architecture which learns to produce a 'code' for an input sequence from which the  
1002 original data can be retrieved. The code is shorter than the original input sequence. Applications include data  
1003 compression and denoising data.

1004 **Back propagation.** A process of utilising the errors from a prediction to update the weights and biases of a neural  
1005 network.

1006 **Batch.** See training batch.



- 1007 **Batch normalisation.** Data normalisation which aligns the means and variances of input data to a model. For  
1008 computational reasons, this is performed separately for each training batch.
- 1009 **Belief state.** The current state of the world which is believed to be true according to a model. A common architecture  
1010 in realtime applications whereby a belief state is updated according to an update function on the basis of new  
1011 observations.
- 1012 **Channel.** An additional dimension to data which is usually not a spatial dimension. Examples include the red, green  
1013 and blue intensity images which comprise a colour image. Another example could be to represent both temperature  
1014 and wind speed as channels.
- 1015 **Classification.** A model which attempts to diagnose or predict the category, label, class or type that an example falls  
1016 within.
- 1017 **Climatology.** Refers to the usual past conditions for a location at a time of year. Usually calculated by temporal mean  
1018 across years of a dataset, for a given time interval within those years (e.g., for a dataset of monthly mean values  
1019 spanning all months of all years from 1990 to 2020, the monthly mean climatology would be obtained by averaging  
1020 across all the Januarys from each year, all the Februarys, etc., to obtain an "average January", an "average February",  
1021 etc.). Climatologies are often used in the same manner as persistence as a baseline prediction against which to measure  
1022 a predictive model. For example, a model predicting a value for January could be compared to the climatological  
1023 monthly mean value for January. This helps answer the question "is my model a better source of information than  
1024 using the average past conditions from this time of year?".
- 1025 **Connectome.** The connections between nodes in a neural network. Examples include fully-connected, partially-  
1026 connected, skip-layer connections, recurrent connections and others. The 'wiring diagram' for the network.
- 1027 **Convolutional neural network.** A deep neural network architecture commonly applied to images which utilises a  
1028 small convolutional (spatially connected) kernel applied in a sliding window fashion.
- 1029 **Data augmentation.** The practice of modifying input data in supervised learning to produce additional examples.  
1030 This can make networks more robust to new inputs and address issues of brittleness to adversarial examples. An  
1031 example of data augmentation is using rotated or reflected versions of the same image as independent training samples.
- 1032 **Data driven.** A generalised term used to indicate a primary reliance or dependence on the collection or analysis of  
1033 data. Used in contrast to process driven or theory driven.
- 1034 **Decision tree.** A tree-like, or flowchart-like, branching model representing a series of decisions and their possible  
1035 consequences. Each internal node represents a 'test' (i.e. decision threshold) and each leaf node represents a class label  
1036 or collection of possible outcomes.
- 1037 **Deep neural network.** A neural network with many layers. Deeper, thinner networks have proven easier to train than  
1038 wider, shallower ones.





- 1039 **DeepONet.** A deep neural network architecture relying on universal approximation theorem to train a neural network  
1040 to represent a mathematical operation (the operator), such as a partial differential equation or dynamic system.
- 1041 **Discriminator model.** A model which distinguishes or discriminates between synthetic data and real-world  
1042 observations. Often used in conjunction with a generator. In this case, the overall goal is to produce a generator which  
1043 is capable of fooling the discriminator, producing highly realistic images. This process is used in Generative  
1044 Adversarial Networks.
- 1045 **Dropout layer.** A neural network layer which is only partially connected, often with a stochastic dropout chance. This  
1046 has been shown experimentally to improve neural network robustness in many architectures by reducing overfitting.
- 1047 **Epoch.** A single complete training pass through all available training data, e.g. learning from all samples, or learning  
1048 from all mini-batches, according to the training strategy. Multiple training epochs will typically be utilised although  
1049 alternative strategies do exist.
- 1050 **Feed-forward network.** A neural network composed of distinct 'layers', where the outputs of one layer never feed  
1051 back into earlier layers. This avoids the needs for any iterative solver approaches and results in a very computationally  
1052 efficient 'forward pass'.
- 1053 **Generative adversarial network.** A two-part neural network architecture comprising a generator and a discriminator,  
1054 which are co-trained to produce realistic outputs which are hard to distinguish from real-world data. The discriminator  
1055 replaces the traditional loss function.
- 1056 **Generator model.** A model which produces a synthetic example of a particular class, such as a synthetic image or  
1057 synthetic language. Examples include language or image generation. These are used as part of Generative Adversarial  
1058 Networks among other applications.
- 1059 **Global receptive field.** Where every part of the input region can influence or stimulate a response in a model (e.g. a  
1060 fully-connected neural network).
- 1061 **GPU.** Graphical Processing Unit. A hardware device specialised for fast matrix operations, originally created to  
1062 support computer graphics, particularly for games.
- 1063 **Gradient boosted decision tree.** Also referred to as extreme gradient boosting. A random forest architecture which  
1064 combines gradient boosting with decision tree ensembles.
- 1065 **Gradient boosting.** An approach to model training where each additional ensemble member attempts to predict the  
1066 cumulative errors of previously trained members.
- 1067 **Graph neural network.** A class of neural networks designed to process data which is described by a graph (or  
1068 tree/network) data structure.
- 1069 **Hidden layer.** A layer which is intermediate between the input layer and the output layer of a network or tree structure.  
1070 Hidden layers may be used to encode 'hidden variables' which are latent to a problem but not able to be directly  
1071 observed, may constitute layers of a deep neural network, or may have other purposes.



- 1072 **Hierarchical temporal aggregation.** A mechanism of composing neural networks which are trained for different lead  
1073 times to produce an optimal prediction at all time horizons.
- 1074 **Hierarchical temporal memory.** Fundamentally different to hierarchical temporal aggregation. A complex deep  
1075 learning architecture which uses time-adjacency pooling.
- 1076 **Hyperparameter.** A parameter which is not derived via training. Examples include the learning rate and the model  
1077 topology.
- 1078 **Hyperparameter search (or Hyperparameter optimisation).** The process of determining optimal hyperparameters.  
1079 This term may also be used to encompass the model selection problem. This process is automated in some cases.
- 1080 **Input layer.** A layer which is composed of input nodes. Typically machine learning models will have one input layer  
1081 at depth zero (i.e. with no preceding layers) and no input nodes at greater depths.
- 1082 **Input node.** A node which represents an input or observed value.
- 1083 **K-fold cross-validation.** A process of changing the validation and test data partitions during different iterations of  
1084 training. This allows more of the training and validation data to be used while minimising overfitting. Some definitions  
1085 include test data in this process but that is not ideal as the final test is no longer statistically independent.
- 1086 **Keras.** A streamlined API for creating neural networks, integrated with Tensorflow. Originally built on the Theano  
1087 framework for general mathematical evaluation. PyTensor and Aesara are related packages.
- 1088 **Kernel trick.** For data sets which are not linearly separable, first multiplying the data by a nonlinear function in a  
1089 higher dimension can result in a linearly separable higher-dimensional data set to which a simpler method can be used  
1090 to model the data.
- 1091 **Knowledge based systems.** A broad term from artificial intelligence meaning a system which that uses reasoning and  
1092 a knowledge base to support decision making. Knowledge is represented explicitly and a reasoning or inference engine  
1093 is used to arrive at new knowledge.
- 1094 **Layer.** In tree or feed-forward network structures (e.g. decision trees and feed-forward neural networks), a layer refers  
1095 to the set of nodes at the same depth within a network.
- 1096 **Leaf node.** Aka output node. A node which does not have any child nodes.
- 1097 **Long short term memory network.** A recurrent neural network architecture which processes sequences of tokens  
1098 utilising a 'memory' component which can store information from tokens early in a sequence for use in prediction of  
1099 tokens much later in a sequence. Typical applications include language prediction and time-series prediction of many  
1100 kinds.
- 1101 **Loss function** (also known as target function, training function, objective function, penalty score, error function,  
1102 heuristic function, minimisation function). A differentiable function which is well-behaved, such that smaller values



1103 represent better model performance and larger values represent worse performance. An example would be the root-  
1104 mean-squared-error of a prediction compared to the truth or target value.

1105 **Mini batch.** A subset or 'mini batch' of the training data. Utilised for multiple reasons, including computational  
1106 efficiency and to reduce overfitting. Aggregate error over a mini-batch is be learned rather than per-sample errors.  
1107 This is the typical contemporary approach. See also training batch for in-depth discussion.

1108 **Neural network.** A composition of 'input nodes', 'connections', 'nodes', 'layers', 'output layers' and 'activation  
1109 functions' which are capable of complex modelling tasks. Originally designed to simulate human neural functioning  
1110 and subsequently applied to a range of applications.

1111 **Node. Aka vertex.** A small data structure in a network, tree or graph structure which is connected by edges. A node  
1112 may represent a real-world value (such as a location) or an abstract value (such as in a neural network), or a decision  
1113 threshold (such as in a decision tree).

1114 **One-hot vector.** A vector of 1s and 0s, in which only one bit is set to 1. Typically produced during the first step in  
1115 machine learning for language processing to create a word or feature embedding in a process called tokenisation or  
1116 encoding. The length of the vector is commonly equal to the number of categories or symbols.

1117 **Output layer.** A layer which comprises the leaf nodes or output nodes of a tree or network.

1118 **Perceptron.** A single-layer neural network architecture for supervised learning of binary classification. Originally  
1119 built as an electronic hardware device encoding weights with potentiometers and learning with motors. A multi-layer  
1120 perceptron is the same thing as an ordinary neural network.

1121 **Persistence.** Refers to the practice of treating some past observation or reanalysis (usually immediately prior to the  
1122 starting point of the prediction period) as the future prediction and "persisting" this one state forward to every  
1123 prediction lead time. The predictive model is then compared to this persistence prediction, essentially assessing the  
1124 performance of the model against a steady state prediction. This, along with climatology, is often used as a baseline  
1125 or bare minimum prediction to beat (i.e., a prediction better than persistence could be considered skilful vs  
1126 persistence). This answers the question "is my model a better source of information than using what happened just  
1127 before now?".

1128 **Physically-informed machine learning. Also known as physics-informed machine learning.** Machine learning is  
1129 considered physically informed when some aspect of physics is included in any way. Examples include adding a  
1130 physical component to the loss function (e.g. to enforce conservation of physical properties) or using an activation  
1131 function with physically realistic properties.

1132 **Predictive step, forward pass, evaluation.** The process of calculating a model prediction from a set of input  
1133 conditions. Distinct from the training phase or back-propagation step.

1134 **PyTorch.** A widely adopted framework for neural networks in Python.



- 1135 **Random forest.** An architecture based on decision tree ensembles where each decision tree is initialised semi-  
1136 randomly and an average of all models is used for prediction. This is typically more accurate than a single decision  
1137 tree but less accurate than a gradient-boosted decision tree and so is now less-used. The term random forest is still  
1138 commonly used when in fact the implementation is a gradient boosted decision tree.
- 1139 **Receptive field.** The size or extent of a region in the input which can influence or stimulate a response in a model,  
1140 e.g. the size of a convolutional kernel, the size of a sliding window
- 1141 **Recurrent network.** A neural network which does pass the output from nodes of the network back into the input of  
1142 others. Infinite recurrence is avoided by setting a specific number of iterations for the recurrence. These are often  
1143 depicted in diagrams as separate layers but the implementation is through internal recurrent connections.
- 1144 **Regression.** A model which attempts to diagnose or predict an exact value by statistically relating example input  
1145 values to desired values.
- 1146 **Relevance vector machine.** A sparse Bayesian model utilising the kernel trick in similar fashion to a support vector  
1147 machine.
- 1148 **Representation error.** Error which is introduced due to the inexactness of representing the real world in the model  
1149 belief state. Examples may include topography smoothing, point-to-grid translations, model grid distortions near the  
1150 poles, or the exclusion of physical characteristics which are not primary to the model.
- 1151 **Residual deep neural network.** A very influential and innovative convolutional deep learning architecture which  
1152 uses a similar concept to gradient boosting. Each layer of the deep network is taken to predict the residual error from  
1153 the previous layers, with skip-connections from earlier layers allowing the training to occur without the issue of  
1154 vanishing gradients.
- 1155 **Sample.** A single training example (e.g. a row of data).
- 1156 **Scale invariance.** A feature of a system, problem or model which means the results and behaviour are the same at any  
1157 scale (e.g., the behaviour does not change if the inputs are multiplied by a common factor).
- 1158 **Scikit-learn.** A popular Python library for machine learning which extends the SciPy framework.
- 1159 **(Stochastic) Gradient descent.** An algorithm by which a neural network is trained using increasingly fine-scale  
1160 adjustments to optimise the accuracy of network prediction. Utilised to find the local minimum of a differentiable  
1161 function.
- 1162 **Supervised learning.** Machine learning is considered 'supervised' when the data is labelled according to a category  
1163 or target value. Classification data have an explicit labelled category. Regression data have an explicit value which is  
1164 being predicted for.
- 1165 **Support vector machine.** A classification model based on finding a hyperplane to separate data utilising the kernel  
1166 trick.



- 1167 **Tensor.** Can be considered as a dense multi-dimensional array or matrix.
- 1168 **Tensorflow.** A widely adopted framework for neural networks in Python.
- 1169 **Test/train/validate split.** Available data is split into three portions. The training data is evaluated and used to update  
1170 model weights. Validation data is evaluated during training and may be used for hyper-parameter search or to guide  
1171 the researcher. Test data is independent (typically well-curated) data used for gold standard evaluation. In reality,  
1172 validation data is sometimes used as test data, but this is not good practice. There are many considerations for  
1173 test/train/validate splitting, such as statistical independence, representation of all classes, and bias. It is important to  
1174 consider what the model is generalising "from" and "to", and ensuring appropriate examples are present in the training  
1175 data and appropriate examples are reserved for validation and test.
- 1176 **Token.** Tokenisation the process of mapping a symbolic or categorical sequence to a numerical representation which  
1177 is suited to a sequence-based machine learning model. Commonly, a vector representation will be utilised for the token  
1178 form. In language processing, either characters or words may be represented as tokens depending on the approach.
- 1179 **TPU.** Tensor Processing Unit. A hardware device specialised for artificial intelligence and machine learning  
1180 applications, in particular neural network operations.
- 1181 **Training batch (or simply batch).** Multiple definitions apply and the use the term has evolved over time. Originally  
1182 used in the context of learning from offline or saved historical data as opposed to online or realtime novel data. In this  
1183 definition, the training batch is the saved data and refers to the whole training set. For example, a robot exploring a  
1184 new environment in real-time must use an online learning technique and could not utilise batch training to map the  
1185 unseen terrain. In more recent use, particularly in the areas of neural network learning, the offline saved data may be  
1186 split into one or more batches (subsets). If one batch (the batch is the entire training set) is used, the aggregate errors  
1187 for the entire training set are used to update the model weights and biases, and the learning algorithm is called batch  
1188 gradient descent. If each example is presented individually, this is called online training (even when historical saved  
1189 data is being used), the weights and biases are updated for from each individual example, and the algorithm used is  
1190 stochastic gradient descent. If the data is divided into multiple batches, this is often referred to equivalently as mini  
1191 batches. The weights and biases are aggregated over each mini batch. This is the most common contemporary  
1192 approach, as it reduces overfitting and is a good balance of training accuracy, avoiding local minima, and  
1193 computational efficiency.
- 1194 **Transfer learning.** The process of training a model first on a related problem, and then conducting further training  
1195 on a more specific problem. Examples could be training a model first in one geographical region and then in another;  
1196 or training first at a low resolution then subsequently at a high resolution. This is frequently done to reduce training  
1197 computation cost for similar problems by re-using the trained weights from a well-performing source model, or to  
1198 overcome a problem of limited data availability by using multiple data sources.



1199 **Transformer network.** A token-sequence architecture which is capable of handling long-range dependencies.  
1200 Initially applied to language processing, it has found effective application in image processing as an alternative to  
1201 convolutional architectures.

1202 **Translation invariance.** A feature of a system, problem or model which means the results and behaviour are the same  
1203 after any spatial translation (i.e., the behaviour does not change if the inputs are shifted spatially to a new location).

1204 **U-Net.** A type of convolutional neural network developed for biomedical image segmentation which has found broad  
1205 application. In the contracting part of the network spatial information is reduced while feature information is increased.  
1206 In the expanding part of the network, feature information is used to inform high-resolution segmentation. The name  
1207 derives from the diagrammatic shape of the network forming a "U".

1208 **Unsupervised learning.** Machine learning is considered 'unsupervised' when data is unlabelled. Examples include  
1209 clustering, association and dimensionality reduction.

1210 **Vanishing Gradient.** At the extremes, nonlinear functions used to calculate gradients can result in gradient values  
1211 which are effectively zero. These small or zero values, once present in the weights and biases of a neural network, can  
1212 entirely suppress information which would in fact be useful, and result in a local minima from which training cannot  
1213 recover. This is particularly relevant to long token-series when long-distance connections are relevant. A variety of  
1214 techniques including alternative activation functions, training weight decay, skip connections and attention  
1215 mechanisms may each or all be utilised to ameliorate this issue.

1216 **Weights and biases.** The parameter values for each neuron which represent the weighting factors to apply to the input  
1217 values, plus an overall bias value for the node.

1218 **XGBoost.** A popular Python library for gradient boosted decision trees.

#### 1219 **Code Availability**

1220 No code was used in the preparation of this review.

#### 1221 **Data Availability**

1222 No data was processed in the preparation of this review.

#### 1223 **Author Contribution**

1224 COdBD researched and wrote Sections 3, 4, 5, 6 and 7, and provided review of sections 2, 8 and 9. TL researched and  
1225 wrote sections 2, 8 and 9, and provided review of sections 3, 4, 5, 6, and 7. COdBD and TL co-wrote sections 1 and  
1226 10.



1227 **Competing Interests**

1228 The authors declare that they have no conflict of interest.

1229 **Acknowledgements**

1230 The authors would like to thank Bethan White, Harrison Cook, Tom Dunstan and Karina Williams for their very  
1231 helpful reviews of early versions of this manuscript.

1232 **References**

1233 Ackmann, J., Düben, P. D., Palmer, T. N., & Smolarkiewicz, P. K. (2020). Machine-learned preconditioners for linear  
1234 solvers in geophysical fluid flows. *arXiv preprint arXiv:2010.02866*.

1235 Adams, S. V., Ford, R. W., Hambley, M., Hobson, J. M., Kavčič, I., Maynard, C. M., ... & Wong, R. (2019). LFRic:  
1236 Meeting the challenges of scalability and performance portability in Weather and Climate models. *Journal of Parallel  
1237 and Distributed Computing*, 132, 383-396.

1238 Arcomano, T., Szunyogh, I., Wikner, A., Pathak, J., Hunt, B. R., & Ott, E. (2022). A Hybrid Approach to Atmospheric  
1239 Modeling That Combines Machine Learning With a Physics-Based Numerical Model. *Journal of Advances in  
1240 Modeling Earth Systems*, 14(3), e2021MS002712.

1241 Atkinson, S. (2020). Bayesian hidden physics models: Uncertainty quantification for discovery of nonlinear partial  
1242 differential operators from data. *arXiv preprint arXiv:2006.04228*.

1243 Bar-Sinai, Y., Hoyer, S., Hickey, J., & Brenner, M. P. (2019). Learning data-driven discretizations for partial  
1244 differential equations. *Proceedings of the National Academy of Sciences*, 116(31), 15344-15349.

1245 Beucler, T., Rasp, S., Pritchard, M., & Gentine, P. (2019). Achieving conservation of energy in neural network  
1246 emulators for climate modeling. *arXiv preprint arXiv:1906.06622*.

1247 Beucler, T., Pritchard, M., Rasp, S., Ott, J., Baldi, P., & Gentine, P. (2021). Enforcing analytic constraints in neural  
1248 networks emulating physical systems. *Physical Review Letters*, 126(9), 098302.

1249 Bhattacharya, K., Hosseini, B., Kovachki, N. B., & Stuart, A. M. (2020). Model reduction and neural networks for  
1250 parametric PDEs. *arXiv preprint arXiv:2005.03180*.

1251 Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X., & Tian, Q. (2022). Pangu-Weather: A 3D High-Resolution Model for  
1252 Fast and Accurate Global Weather Forecast. *arXiv preprint arXiv:2211.02556*.

1253 Bihlo, A., & Popovych, R. O. (2022). Physics-informed neural networks for the shallow-water equations on the sphere.  
1254 *Journal of Computational Physics*, 456, 111024.

1255 Bolton, T., & Zanna, L. (2019). Applications of deep learning to ocean data inference and subgrid parameterization.  
1256 *Journal of Advances in Modeling Earth Systems*, 11(1), 376-399.



- 1257 Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32.
- 1258 Brenowitz, N. D., & Bretherton, C. S. (2018). Prognostic validation of a neural network unified physics  
1259 parameterization. *Geophysical Research Letters*, 45(12), 6289-6298.
- 1260 Brenowitz, N. D., & Bretherton, C. S. (2019). Spatially extended tests of a neural network parametrization trained by  
1261 coarse-graining. *Journal of Advances in Modeling Earth Systems*, 11(8), 2728-2744.
- 1262 Brenowitz, N. D., Beucler, T., Pritchard, M., & Bretherton, C. S. (2020). Interpreting and stabilizing machine-learning  
1263 parametrizations of convection. *Journal of the Atmospheric Sciences*, 77(12), 4357-4375.
- 1264 Brenowitz, N. D., Henn, B., McGibbon, J., Clark, S. K., Kwa, A., Perkins, W. A., ... & Bretherton, C. S. (2020).  
1265 Machine learning climate model dynamics: Offline versus online performance. *arXiv preprint arXiv:2011.03081*.
- 1266 Brenowitz, N. D., Perkins, W. A., Nugent, J. M., Watt-Meyer, O., Clark, S. K., Kwa, A., ... & Bretherton, C. S. (2022).  
1267 Emulating Fast Processes in Climate Models. *arXiv preprint arXiv:2211.10774*.
- 1268 Chantry, M., Christensen, H., Dueben, P., & Palmer, T. (2021). Opportunities and challenges for machine learning in  
1269 weather and climate modelling: hard, medium and soft AI. *Philosophical Transactions of the Royal Society A*,  
1270 379(2194), 20200083.
- 1271 Chantry, M., Hatfield, S., Dueben, P., Polichtchouk, I., & Palmer, T. (2021). Machine learning emulation of gravity  
1272 wave drag in numerical weather forecasting. *Journal of Advances in Modeling Earth Systems*, 13(7), e2021MS002477.
- 1273 Chase, R. J., Harrison, D. R., Burke, A., Lackmann, G. M., & McGovern, A. (2022a). A Machine Learning Tutorial  
1274 for Operational Meteorology, Part I: Traditional Machine Learning. *arXiv preprint arXiv:2204.07492*.
- 1275 Chase, R. J., Harrison, D. R., Lackmann, G., & McGovern, A. (2022b). A Machine Learning Tutorial for Operational  
1276 Meteorology, Part II: Neural Networks and Deep Learning. *arXiv preprint arXiv:2211.00147*.
- 1277 Chattopadhyay, A., Subel, A., & Hassanzadeh, P. (2020). Data-driven super-parameterization using deep learning:  
1278 Experimentation with multiscale Lorenz 96 systems and transfer learning. *Journal of Advances in Modeling Earth  
1279 Systems*, 12(11), e2020MS002084.
- 1280 Chevallier, F., Ch eruy, F., Scott, N. A., & Ch edin, A. (1998). A neural network approach for a fast and accurate  
1281 computation of a longwave radiative budget. *Journal of applied meteorology*, 37(11), 1385-1397.
- 1282 Clare, M. C., Jamil, O., & Morcrette, C. J. (2021). Combining distribution-based neural networks to predict weather  
1283 forecast probabilities. *Quarterly Journal of the Royal Meteorological Society*, 147(741), 4337-4357.
- 1284 Cuomo, S., Di Cola, V. S., Giampaolo, F., Rozza, G., Raissi, M., & Piccialli, F. (2022). Scientific Machine Learning  
1285 through Physics-Informed Neural Networks: Where we are and What's next. *arXiv preprint arXiv:2201.05624*.
- 1286 De B ezenac, E., Pajot, A., & Gallinari, P. (2017). *Towards a hybrid approach to physical process modeling*. Technical  
1287 report.





- 1288 Deng, J., Dong, W., Socher, R., Li, L. J., Li, K., & Fei-Fei, L. (2009, June). Imagenet: A large-scale hierarchical image  
1289 database. In *2009 IEEE conference on computer vision and pattern recognition* (pp. 248-255). Ieee.
- 1290 Dijkstra, H. A., Petersik, P., Hernández-García, E., & López, C. (2019). The application of machine learning  
1291 techniques to improve El Niño prediction skill. *Frontiers in Physics*, 153.
- 1292 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. (2020). An  
1293 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- 1294 Dueben, P. D., & Bauer, P. (2018). Challenges and design choices for global weather and climate models based on  
1295 machine learning. *Geoscientific Model Development*, 11(10), 3999-4009.
- 1296 Dueben, P. D., Schultz, M. G., Chantry, M., Gagne, D. J., Hall, D. M., & McGovern, A. (2022). Challenges and  
1297 Benchmark Datasets for Machine Learning in the Atmospheric Sciences: Definition, Status, and Outlook. *Artificial  
1298 Intelligence for the Earth Systems*, 1(3), e210002.
- 1299 ECMWF. (2018). Ifs documentation (cy45r1). Retrieved from [https://www.ecmwf.int/en/publications/ifs-  
1300 documentation](https://www.ecmwf.int/en/publications/ifs-<br/>1300 documentation) accessed 7<sup>th</sup> February 2023
- 1301 Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., & Taylor, K. E. (2016). Overview of  
1302 the Coupled Model Intercomparison Project Phase 6 (CMIP6) experimental design and organization. *Geoscientific  
1303 Model Development*, 9(5), 1937-1958.
- 1304 Flora, M., Potvin, C., McGovern, A., & Handler, S. (2022). Comparing Explanation Methods for Traditional Machine  
1305 Learning Models Part 2: Quantifying Model Explainability Faithfulness and Improvements with Dimensionality  
1306 Reduction. *arXiv preprint arXiv:2211.10378*.
- 1307 Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of statistics*, 1189-1232.
- 1308 Fuhg, J. N., Kalogeris, I., Fau, A., & Bouklas, N. (2022). Interval and fuzzy physics-informed neural networks for  
1309 uncertain fields. *Probabilistic Engineering Mechanics*, 68, 103240.
- 1310 Gagne, D. J., McCandless, T., Kosovic, B., DeCastro, A., Loft, R., Haupt, S. E., & Yang, B. (2019, December).  
1311 Machine learning parameterization of the surface layer: bridging the observation-modeling gap. In *AGU Fall Meeting  
1312 Abstracts* (Vol. 2019, pp. IN44A-04).
- 1313 Gagne, D. J., Christensen, H. M., Subramanian, A. C., & Monahan, A. H. (2020). Machine learning for stochastic  
1314 parameterization: Generative adversarial networks in the Lorenz'96 model. *Journal of Advances in Modeling Earth  
1315 Systems*, 12(3), e2019MS001896.
- 1316 Gagne, D. J., Chen, C. C., & Gettelman, A. (2020, January). Emulation of bin Microphysical Processes with machine  
1317 learning. In *100th American Meteorological Society Annual Meeting*. AMS.
- 1318 George, T., Gupta, A., & Sarin, V. (2008, December). A recommendation system for preconditioned iterative solvers.  
1319 In *2008 Eighth IEEE International Conference on Data Mining* (pp. 803-808). IEEE.



- 1320 Guttelman, A., Gagne, D. J., Chen, C. C., Christensen, M. W., Lebo, Z. J., Morrison, H., & Gantos, G. (2021). Machine  
1321 learning the warm rain process. *Journal of Advances in Modeling Earth Systems*, 13(2), e2020MS002268.
- 1322 Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., ... & Bengio, Y. (2020). Generative  
1323 adversarial networks. *Communications of the ACM*, 63(11), 139-144.
- 1324 Goodfellow, I., Yoshua B., & Aaron C. (2016). Deep learning. *MIT press*.
- 1325 Grigo, C., & Koutsourelakis, P. S. (2019). A physics-aware, probabilistic machine learning framework for coarse-  
1326 graining high-dimensional systems in the Small Data regime. *Journal of Computational Physics*, 397, 108842.
- 1327 Ham, Y. G., Kim, J. H., & Luo, J. J. (2019). Deep learning for multi-year ENSO forecasts. *Nature*, 573(7775), 568-  
1328 572.
- 1329 Ham, Y. G., Kim, J. H., Kim, E. S., & On, K. W. (2021). Unified deep learning model for El Niño/Southern Oscillation  
1330 forecasts by incorporating seasonality in climate data. *Science Bulletin*, 66(13), 1358-1366.
- 1331 Han, Y., Zhang, G. J., Huang, X., & Wang, Y. (2020). A moist physics parameterization based on deep learning.  
1332 *Journal of Advances in Modeling Earth Systems*, 12(9), e2020MS002076.
- 1333 Harder, P., Watson-Parris, D., Stier, P., Strassel, D., Gauger, N. R., & Keuper, J. (2022). Physics-informed learning  
1334 of aerosol microphysics. *Environmental Data Science*, 1, e20.
- 1335 Harris, L., Chen, X., Putman, W., Zhou, L., & Chen, J. H. (2021). A scientific description of the GFDL finite-volume  
1336 cubed-sphere dynamical core.
- 1337 Hastie, T., Tibshirani, R., Friedman, J. H., & Friedman, J. H. (2009). *The elements of statistical learning: data mining,*  
1338 *inference, and prediction* (Vol. 2, pp. 1-758). New York: springer.
- 1339 He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the*  
1340 *IEEE conference on computer vision and pattern recognition* (pp. 770-778).
- 1341 Hewamalage, H., Ackermann, K., & Bergmeir, C. (2022). Forecast Evaluation for Data Scientists: Common Pitfalls  
1342 and Best Practices. *arXiv preprint arXiv:2203.10716*.
- 1343 Hochreiter, S., & Schmidhuber, J. (1996). LSTM can solve hard long time lag problems. *Advances in neural*  
1344 *information processing systems*, 9.
- 1345 Holloway, A., & Chen, T. Y. (2007, May). Neural networks for predicting the behavior of preconditioned iterative  
1346 solvers. In *International Conference on Computational Science* (pp. 302-309). Springer, Berlin, Heidelberg.
- 1347 Hornik, K., Stinchcombe, M., & White, H. (1989). Multilayer feedforward networks are universal  
1348 approximators. *Neural networks*, 2(5), 359-366.
- 1349 Huang, Z., England, M., Davenport, J. H., & Paulson, L. C. (2016, September). Using machine learning to decide  
1350 when to precondition cylindrical algebraic decomposition with Groebner bases. In *2016 18th International Symposium*  
1351 *on Symbolic and Numeric Algorithms for Scientific Computing (SYNASC)* (pp. 45-52). IEEE.



- 1352 Karpatne, A., Atluri, G., Faghmous, J. H., Steinbach, M., Banerjee, A., Ganguly, A., ... & Kumar, V. (2017). Theory-  
1353 guided data science: A new paradigm for scientific discovery from data. *IEEE Transactions on knowledge and data*  
1354 *engineering*, 29(10), 2318-2331.
- 1355 Karunasinghe, D. S., & Liong, S. Y. (2006). Chaotic time series prediction with a global model: Artificial neural  
1356 network. *Journal of Hydrology*, 323(1-4), 92-105.
- 1357 Kashinath, K., Mustafa, M., Albert, A., Wu, J. L., Jiang, C., Esmaeilzadeh, S., ... & Prabhat. (2021). Physics-informed  
1358 machine learning: case studies for weather and climate modelling. *Philosophical Transactions of the Royal Society A*,  
1359 379(2194), 20200093.
- 1360 Keisler, R. (2022). Forecasting Global Weather with Graph Neural Networks. *arXiv preprint arXiv:2202.07575*.
- 1361 Kim, J., Kwon, M., Kim, S. D., Kug, J. S., Ryu, J. G., & Kim, J. (2022). Spatiotemporal neural network with attention  
1362 mechanism for El Niño forecasts. *Scientific Reports*, 12(1), 1-15.
- 1363 Kochkov, D., Smith, J. A., Alieva, A., Wang, Q., Brenner, M. P., & Hoyer, S. (2021). Machine learning–accelerated  
1364 computational fluid dynamics. *Proceedings of the National Academy of Sciences*, 118(21), e2101784118.
- 1365 Krasnopolsky, V. M., Chalikov, D. V., & Tolman, H. L. (2002). A neural network technique to improve computational  
1366 efficiency of numerical oceanic models. *Ocean Modelling*, 4(3-4), 363-383.
- 1367 Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Chalikov, D. V. (2005). New approach to calculation of atmospheric  
1368 model physics: Accurate and fast neural network emulation of longwave radiation in a climate model. *Monthly*  
1369 *Weather Review*, 133(5), 1370-1383.
- 1370 Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Belochitski, A. A. (2013). Using ensemble of neural networks to learn  
1371 stochastic convection parameterizations for climate and numerical weather prediction models from data simulated by  
1372 a cloud resolving model. *Advances in Artificial Neural Systems*, 2013.
- 1373 Kuefler, E., & Chen, T. Y. (2008, June). On using reinforcement learning to solve sparse linear systems. In  
1374 *International Conference on Computational Science* (pp. 955-964). Springer, Berlin, Heidelberg.
- 1375 Ladický, L. U., Jeong, S., Solenthaler, B., Pollefeys, M., & Gross, M. (2015). Data-driven fluid simulations using  
1376 regression forests. *ACM Transactions on Graphics (TOG)*, 34(6), 1-9.
- 1377 Lam, R., Sanchez-Gonzalez, A., Willson, M., Wirsberger, P., Fortunato, M., Pritzel, A., ... & Battaglia, P. (2022).  
1378 GraphCast: Learning skillful medium-range global weather forecasting. *arXiv preprint arXiv:2212.12794*.
- 1379 Lanthaler, S., Mishra, S., & Karniadakis, G. E. (2022). Error estimates for deeponets: A deep learning framework in  
1380 infinite dimensions. *Transactions of Mathematics and Its Applications*, 6(1), tnac001.
- 1381 LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition.  
1382 *Proceedings of the IEEE*, 86(11), 2278-2324.



- 1383 Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Stuart, A., Bhattacharya, K., & Anandkumar, A. (2020a). Multipole  
1384 graph neural operator for parametric partial differential equations. *Advances in Neural Information Processing*  
1385 *Systems*, 33, 6755-6766.
- 1386 Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Bhattacharya, K., Stuart, A., & Anandkumar, A. (2020b). Neural  
1387 operator: Graph kernel network for partial differential equations. *arXiv preprint arXiv:2003.03485*.
- 1388 Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Bhattacharya, K., Stuart, A., & Anandkumar, A. (2020c). Fourier  
1389 neural operator for parametric partial differential equations. *arXiv preprint arXiv:2010.08895*.
- 1390 Lorenz, E. N. (1996). Predictability: A problem partly solved, in Proceedings of Seminar on Predictability, 4–8  
1391 September 1995.
- 1392 Lu, L., Jin, P., & Karniadakis, G. E. (2019). Deeponet: Learning nonlinear operators for identifying differential  
1393 equations based on the universal approximation theorem of operators. *arXiv preprint arXiv:1910.03193*.
- 1394 McCulloch, W. S., & Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *The bulletin of*  
1395 *mathematical biophysics*, 5(4), 115-133.
- 1396 Meyer, D., Hogan, R. J., Dueben, P. D., & Mason, S. L. (2022). Machine learning emulation of 3D cloud radiative  
1397 effects. *Journal of Advances in Modeling Earth Systems*, 14(3), e2021MS002550.
- 1398 Mooers, G., Pritchard, M., Beucler, T., Ott, J., Yacalis, G., Baldi, P., & Gentine, P. (2021). Assessing the Potential of  
1399 Deep Learning for Emulating Cloud Superparameterization in Climate Models With Real-Geography Boundary  
1400 Conditions. *Journal of Advances in Modeling Earth Systems*, 13(5), e2020MS002385.
- 1401 Mudigonda, M., Kim, S., Mahesh, A., Kahou, S., Kashinath, K., Williams, D., ... & Prabhat, M. (2017). Segmenting  
1402 and tracking extreme climate events using neural networks. In *Deep Learning for Physical Sciences (DLPS)*  
1403 *Workshop, held with NIPS Conference*.
- 1404 Nelsen, N. H., & Stuart, A. M. (2021). The random feature model for input-output maps between banach spaces. *SIAM*  
1405 *Journal on Scientific Computing*, 43(5), A3212-A3243.
- 1406 Nguyen, T., Brandstetter, J., Kapoor, A., Gupta, J. K., & Grover, A. (2023). ClimaX: A foundation model for weather  
1407 and climate. *arXiv preprint arXiv:2301.10343*. Nielsen, A. H., Iosifidis, A., & Karstoft, H. (2022). Forecasting large-  
1408 scale circulation regimes using deformable convolutional neural networks and global spatiotemporal climate data.  
1409 *Scientific Reports*, 12(1), 1-12.
- 1410 O’Gorman, P. A., & Dwyer, J. G. (2018). Using machine learning to parameterize moist convection: Potential for  
1411 modeling of climate, climate change, and extreme events. *Journal of Advances in Modeling Earth Systems*, 10(10),  
1412 2548-2563.
- 1413 O’Leary, J., Paulson, J. A., & Mesbah, A. (2022). Stochastic physics-informed neural ordinary differential equations.  
1414 *Journal of Computational Physics*, 468, 111466.



- 1415 Ott, J., Pritchard, M., Best, N., Linstead, E., Curcic, M., & Baldi, P. (2020). A Fortran-Keras deep learning bridge for  
1416 scientific computing. *Scientific Programming*, 2020.
- 1417 Palmer, T. (2020). A vision for numerical weather prediction in 2030. *arXiv preprint arXiv:2007.04830*.
- 1418 Patel, R. G., Trask, N. A., Wood, M. A., & Cyr, E. C. (2021). A physics-informed operator regression framework for  
1419 extracting data-driven continuum models. *Computer Methods in Applied Mechanics and Engineering*, 373, 113500.
- 1420 Pathak, J., Subramanian, S., Harrington, P., Raja, S., Chattopadhyay, A., Mardani, M., ... & Anandkumar, A. (2022).  
1421 Fourcastnet: A global data-driven high-resolution weather model using adaptive fourier neural operators. *arXiv  
1422 preprint arXiv:2202.11214*.
- 1423 Peairs, L., & Chen, T. Y. (2011). Using reinforcement learning to vary the m in GMRES (m). *Procedia Computer  
1424 Science*, 4, 2257-2266.
- 1425 Pincus, R., Mlawer, E. J., & Delamere, J. S. (2019). Balancing accuracy, efficiency, and flexibility in radiation  
1426 calculations for dynamical models. *Journal of Advances in Modeling Earth Systems*, 11(10), 3074-3089.
- 1427 Psaros, A. F., Meng, X., Zou, Z., Guo, L., & Karniadakis, G. E. (2022). Uncertainty quantification in scientific  
1428 machine learning: Methods, metrics, and comparisons. *arXiv preprint arXiv:2201.07766*.
- 1429 Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid processes in climate models.  
1430 *Proceedings of the National Academy of Sciences*, 115(39), 9684-9689.
- 1431 Rasp, S. (2020). Coupled online learning as a way to tackle instabilities and biases in neural network  
1432 parameterizations: general algorithms and Lorenz 96 case study (v1. 0). *Geoscientific Model Development*, 13(5),  
1433 2185-2196.
- 1434 Rasp, S., Dueben, P. D., Scher, S., Weyn, J. A., Mouatadid, S., & Thuerey, N. (2020). WeatherBench: a benchmark  
1435 data set for data-driven weather forecasting. *Journal of Advances in Modeling Earth Systems*, 12(11),  
1436 e2020MS002203.
- 1437 Rasp, S., & Thuerey, N. (2021). Data-driven medium-range weather prediction with a resnet pretrained on climate  
1438 simulations: A new model for weatherbench. *Journal of Advances in Modeling Earth Systems*, 13(2),  
1439 e2020MS002405.
- 1440 Ravuri, S., Lenc, K., Willson, M., Kangin, D., Lam, R., Mirowski, P., ... & Mohamed, S. (2021). Skilful precipitation  
1441 nowcasting using deep generative models of radar. *Nature*, 597(7878), 672-677.
- 1442 Rizzuti, G., Siahkoobi, A., & Herrmann, F. J. (2019, June). Learned iterative solvers for the Helmholtz equation. In  
1443 *81st EAGE Conference and Exhibition 2019* (Vol. 2019, No. 1, pp. 1-5). European Association of Geoscientists &  
1444 Engineers.



- 1445 Ronneberger, O., Fischer, P., & Brox, T. (2015, October). U-net: Convolutional networks for biomedical image  
1446 segmentation. In *International Conference on Medical image computing and computer-assisted intervention* (pp. 234-  
1447 241). Springer, Cham.
- 1448 Russell S. & Norvig P (2021). *Artificial Intelligence: A Modern Approach* (Fourth Global Edition). Pearson Education.
- 1449 Scher, S. (2018). Toward data-driven weather and climate forecasting: Approximating a simple general circulation  
1450 model with deep learning. *Geophysical Research Letters*, 45(22), 12-616.
- 1451 Scher, S., & Messori, G. (2019). Weather and climate forecasting with neural networks: using general circulation  
1452 models (GCMs) with different complexity as a study ground. *Geoscientific Model Development*, 12(7), 2797-2809.
- 1453 Sønderby, C. K., Espeholt, L., Heek, J., Dehghani, M., Oliver, A., Salimans, T., ... & Kalchbrenner, N. (2020). Metnet:  
1454 A neural weather model for precipitation forecasting. *arXiv preprint arXiv:2003.12140*.
- 1455 Taylor, J., & Feng, M. (2022). A Deep Learning Model for Forecasting Global Monthly Mean Sea Surface  
1456 Temperature Anomalies. *arXiv preprint arXiv:2202.09967*.
- 1457 Tibshirani, R., & Friedman, J. H. (2001). *The elements of statistical learning [electronic resource]: data mining,  
1458 inference, and prediction: with 200 full-color illustrations*. Springer.
- 1459 Tompson, J., Schlachter, K., Sprechmann, P., & Perlin, K. (2017, July). Accelerating eulerian fluid simulation with  
1460 convolutional networks. In *International Conference on Machine Learning* (pp. 3424-3433). PMLR.
- 1461 Turing, A. M., Computing Machinery and Intelligence, *Mind*, Volume LIX, Issue 236, October 1950, Pages 433–460,  
1462 <https://doi.org/10.1093/mind/LIX.236.433>
- 1463 Ukkonen, P., & Mäkelä, A. (2019). Evaluation of machine learning classifiers for predicting deep convection. *Journal  
1464 of Advances in Modeling Earth Systems*, 11(6), 1784-1802.
- 1465 Ukkonen, P., Pincus, R., Hogan, R. J., Pagh Nielsen, K., & Kaas, E. (2020). Accelerating radiation computations for  
1466 dynamical models with targeted machine learning and code optimization. *Journal of Advances in Modeling Earth  
1467 Systems*, 12(12), e2020MS002226.
- 1468 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is  
1469 all you need. *Advances in neural information processing systems*, 30.
- 1470 Vlachas, P. R., Byeon, W., Wan, Z. Y., Sapsis, T. P., & Koumoutsakos, P. (2018). Data-driven forecasting of high-  
1471 dimensional chaotic systems with long short-term memory networks. *Proceedings of the Royal Society A:  
1472 Mathematical, Physical and Engineering Sciences*, 474(2213), 20170844.
- 1473 Walters, D., Boutle, I., Brooks, M., Melvin, T., Stratton, R., Vosper, S., ... & Xavier, P. (2017). The Met Office unified  
1474 model global atmosphere 6.0/6.1 and JULES global land 6.0/6.1 configurations. *Geoscientific Model Development*,  
1475 10(4), 1487-1520.



- 1476 Wang, S., Wang, H., & Perdikaris, P. (2021). Learning the solution operator of parametric partial differential equations  
1477 with physics-informed DeepONets. *Science advances*, 7(40), eabi8605.
- 1478 Wang, X., Han, Y., Xue, W., Yang, G., & Zhang, G. J. (2022a). Stable climate simulations using a realistic general  
1479 circulation model with neural network parameterizations for atmospheric moist physics and radiation processes.  
1480 *Geoscientific Model Development*, 15(9), 3923-3940.
- 1481 Wang, S., Sankaran, S., & Perdikaris, P. (2022b). Respecting causality is all you need for training physics-informed  
1482 neural networks. *arXiv preprint arXiv:2203.07404*.
- 1483 Watson, P. A. (2022). Machine learning applications for weather and climate need greater focus on extremes.  
1484 *Environmental Research Letters*, 17(11), 111004.
- 1485 Watt-Meyer, O., Brenowitz, N. D., Clark, S. K., Henn, B., Kwa, A., McGibbon, J., ... & Bretherton, C. S. (2021).  
1486 Correcting weather and climate models by machine learning nudged historical simulations. *Geophysical Research*  
1487 *Letters*, 48(15), e2021GL092555.
- 1488 Werbos, P. (1974). Beyond regression: new tools for prediction and analysis in the behavioral sciences. *Ph. D.*  
1489 *dissertation, Harvard University*.
- 1490 Werbos, P. J. (1990). Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10),  
1491 1550-1560.
- 1492 Weyn, J. A., Durran, D. R., & Caruana, R. (2019). Can machines learn to predict weather? Using deep learning to  
1493 predict gridded 500-hPa geopotential height from historical weather data. *Journal of Advances in Modeling Earth*  
1494 *Systems*, 11(8), 2680-2693.
- 1495 Weyn, J. A., Durran, D. R., & Caruana, R. (2020). Improving data-driven global weather prediction using deep  
1496 convolutional neural networks on a cubed sphere. *Journal of Advances in Modeling Earth Systems*, 12(9),  
1497 e2020MS002109.
- 1498 Weyn, J. A., Durran, D. R., Caruana, R., & Cresswell-Clay, N. (2021). Sub-seasonal forecasting with a large ensemble  
1499 of deep-learning weather prediction models. *Journal of Advances in Modeling Earth Systems*, 13(7)
- 1500 Wikner, A., Pathak, J., Hunt, B., Girvan, M., Arcomano, T., Szunyogh, I., ... & Ott, E. (2020). Combining machine  
1501 learning with knowledge-based modeling for scalable forecasting and subgrid-scale closure of large, complex,  
1502 spatiotemporal systems. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 30(5), 053111.
- 1503 Wu, K., & Xiu, D. (2020). Data-driven deep learning of partial differential equations in modal space. *Journal of*  
1504 *Computational Physics*, 408, 109307.
- 1505 Yamada, K., Katagiri, T., Takizawa, H., Minami, K., Yokokawa, M., Nagai, T., & Ogino, M. (2018, November).  
1506 Preconditioner auto-tuning using deep learning for sparse iterative algorithms. In *2018 Sixth International Symposium*  
1507 *on Computing and Networking Workshops (CANDARW)* (pp. 257-262). IEEE.



- 1508 Yang, C., Yang, X., & Xiao, X. (2016). Data-driven projection method in fluid simulation. *Computer Animation and*  
1509 *Virtual Worlds*, 27(3-4), 415-424.
- 1510 Yeo, K., Grullon, D. E., Sun, F. K., Boning, D. S., & Kalagnanam, J. R. (2021). Variational inference formulation for  
1511 a model-free simulation of a dynamical system with unknown parameters by a recurrent neural network. *SIAM Journal*  
1512 *on Scientific Computing*, 43(2), A1305-A1335.
- 1513 Yuval, J., & O’Gorman, P. A. (2020). Stable machine-learning parameterization of subgrid processes for climate  
1514 modeling at a range of resolutions. *Nature communications*, 11(1), 1-10.
- 1515 Yuval, J., O’Gorman, P. A., & Hill, C. N. (2021). Use of neural networks for stable, accurate and physically consistent  
1516 parameterization of subgrid atmospheric processes with good performance at reduced precision. *Geophysical*  
1517 *Research Letters*, 48(6), e2020GL091363.
- 1518 Zanna, L., & Bolton, T. (2020). Data-driven equation discovery of ocean mesoscale closures. *Geophysical Research*  
1519 *Letters*, 47(17), e2020GL088376.
- 1520 Zhong, X., Ma, Z., Yao, Y., Xu, L., Wu, Y., & Wang, Z. (2023). WRF–ML v1. 0: a bridge between WRF v4. 3 and  
1521 machine learning parameterizations and its application to atmospheric radiative transfer. *Geoscientific Model*  
1522 *Development*, 16(1), 199-209.
- 1523 Zhou, L., Lin, S. J., Chen, J. H., Harris, L. M., Chen, X., & Rees, S. L. (2019). Toward convective-scale prediction  
1524 within the next generation global prediction system. *Bulletin of the American Meteorological Society*, 100(7), 1225-  
1525 1243.
- 1526