

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

3 Catherine O. de Burgh-Day¹ & Tennessee Leeuwenburg¹

4 ¹The Bureau of Meteorology, 700 Collins St Docklands, Victoria, Australia

5 *Correspondence to:* Catherine O. de Burgh-Day (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 parameterization scheme emulation and replacement, and more recently
9 even to full ML-based weather and climate prediction models. While ML has been used in this space for more than
10 25 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 terms, methodologies, and ethical
14 considerations. Finally, we discuss some potentially beneficial future research directions. Our aim is to provide a
15 primer for researchers and model developers to rapidly familiarize and update themselves with the world of ML in the
16 context of weather and climate models.

17 **1. Introduction**

18 Current state-of-the-art weather and climate models use numerical methods to solve equations representing the
19 dynamics of the atmosphere and ocean on meshed grids. The grid-scale effects of processes that are too small to be
20 resolved are either represented by parametrization schemes or are prescribed. These numerical weather and climate
21 forecasts are computationally costly and are not easy to implement on specialized compute resources such as GPUs
22 (although there are efforts underway to do so, for example in LFRic (Adams et al. 2019)). One of the main approaches
23 to improving forecast accuracy is to increase model resolution (reduced timestep between model increments and/or
24 decreased grid spacing), but due to the high computational cost of this approach, improvements in model skill are
25 hampered by the finite supercomputer capacity available. An additional pathway to improve skill is to improve the
26 understanding and representation of subgrid-scale processes, however this is again a potentially computationally costly
27 exercise.

28 In the remainder of this introduction, we overview the state of machine learning in weather and climate research
29 without always providing references; we instead provide relevant references in the detailed sections that follow.

30 Machine learning is an increasingly powerful and popular tool. It has proven to be computationally efficient, as well
31 as being an accurate way to model subgrid-scale processes. The term “Machine learning” (ML) was first coined by

32 Arthur Samuel in 1952 to refer to a “field of study that gives computers the ability to learn without being explicitly
33 programmed”¹. Learning by example is the defining characteristic of ML.

34 The growing potential for ML in weather and climate modelling is being increasingly recognized by meteorological
35 agencies and researchers around the world. The former is evidenced by the development of strategies and frameworks
36 to better support the development of ML research, such as the Data Science Framework recently published by the Met
37 Office in the UK². The latter is made clear by the explosion in publications from academia, government agencies and
38 private industry in this space, as demonstrated by the rest of this review. Figure 1 shows the number of publications
39 cited in this review using different categories of ML algorithms by year, and clearly illustrates the increase in the
40 uptake of ML methods by the research community.

41 This is not necessarily an unbiased sample of the use of different architectures in the literature, since the selection of
42 papers cited in this review focuses on telling the story of the growth of the use of ML in weather and climate modelling
43 over time, rather than being a comprehensive list of all uses of ML in the literature.

44 There are established techniques and aspects of the weather and climate modelling lifecycle that would already be
45 considered ML by many. For example, linear regression^{†3}, principal component analysis, correlations, and the
46 calculation of teleconnections can all be considered types of ML. Data Assimilation techniques could also be
47 considered a form of ML. There are, however, other classes of ML (e.g. Neural Networks[†], Decision Trees[†], etc.)
48 which are much less widely used within the weather and climate modelling space and have great potential to be of
49 benefit. There is growing interest in, and increasingly effective application of, these ML techniques to take the place
50 of more traditional approaches to modelling. The potential for ML in weather and climate modelling extends all the
51 way from replacement of individual sub-components of the model (to improve accuracy and reduce computational
52 cost) to full replacement of the entire numerical model.

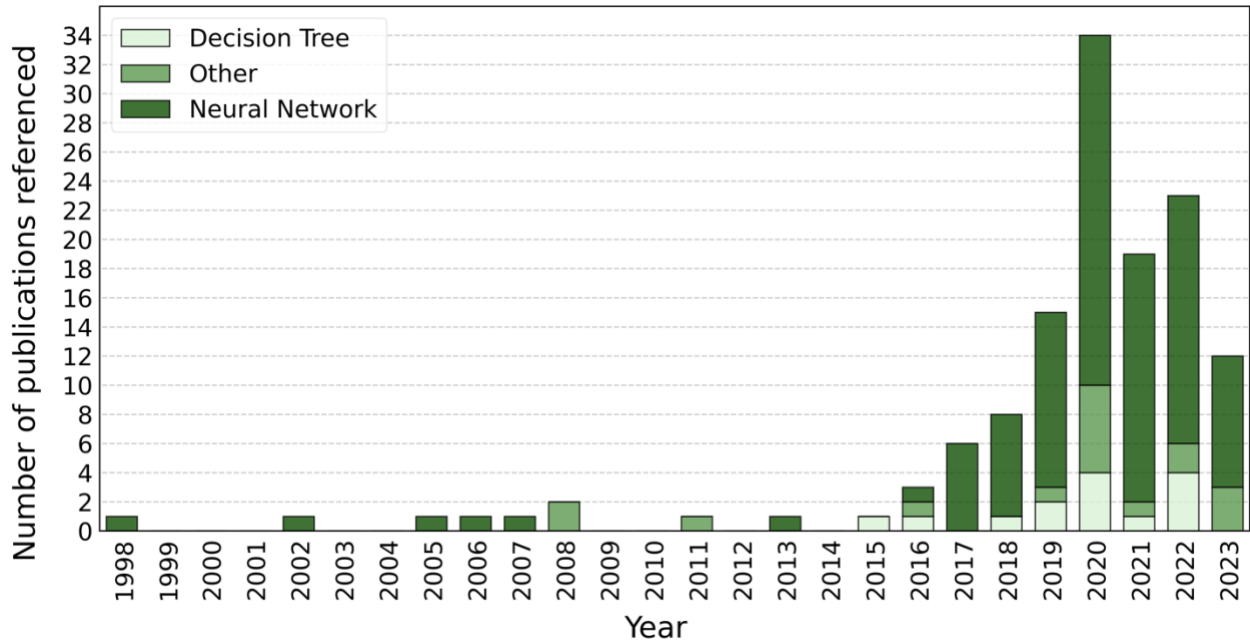
53 While ML models are typically computationally costly during training, they can provide very fast predictions at
54 inference[†] time, especially on GPU hardware. They often also avoid the need to have full understanding of the
55 processes being represented and can learn and infer complex relationships without any need for them to be explicitly
56 encoded. These properties make ML an attractive alternative to traditional parametrization, numerical solver, and
57 modelling methods.

58 Neural Networks (NNs, explained further in Section 2.1) in particular are an increasingly favored alternative approach
59 for representing sub-grid-scale processes or replacing numerical models entirely. They consist of several
60 interconnected layers of nonlinear nodes[†], with the number of intermediate layers depending on the complexity of the
61 system being represented. These nodes allow for the encoding of an arbitrary number of interrelationships between
62 arbitrary parameters to represent the system, removing the need to explicitly encode these interrelationships into a
63 parameterization or numerical model.

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

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

³ Henceforth, the first occurrence of each term described in the glossary is marked with the symbol "†"



64
 65 Figure 1: A stacked bar graph of the number of publications cited in this review using different categories of ML algorithms by per
 66 year. For a description of Neural Networks and Decision Trees see Section 2.1 and 2.2 respectively. The ‘Other’ category is a
 67 collection of ML model types other than decision trees and neural networks, each of which only had one or two examples of use in
 68 this review. This included custom supervised and self-supervised algorithms, support vector machines and relevance vector
 69 machine models, regression models, unsupervised learning models, reservoir computing models and non-NN gaussian models.
 70 This figure includes all references from this review except for: seminal ML papers that are on new ML methods (e.g., foundational
 71 ML papers from outside the domain of weather and climate modelling), review papers, any paper cited that concerns a topic which
 72 is out of scope (e.g., nowcasting), and any other paper which does not present a new method directly applicable to weather and
 73 climate modelling. The full table of citations is provided in the appendix.

74
 75
 76 One challenge that must be overcome before there will be more widespread acceptance of ML as an alternative to
 77 traditional modelling methods is that ML is seen as lacking interpretability. Most ML models do not explicitly
 78 represent the physical processes they are simulating, although physics constrained ML is a new and growing field
 79 which goes some way to addressing this (see Section 6). Furthermore, the techniques available to gain insight into the
 80 relative importance and predictive mechanism of each predictor (i.e. the model outputs) are limited. In contrast,
 81 traditional models are usually driven by some understanding and/or representation of the physical mechanisms and
 82 processes which are occurring. This makes it possible to more easily gain insight into what physical drivers could
 83 explain a given output. The “black box” nature of many current ML approaches to parametrization makes them an
 84 unpopular choice for many researchers (and can be off-putting for decision makers) since, for example, explaining
 85 what went wrong in a model after a bad forecast can be more challenging if there are processes in the model which
 86 are not, and cannot, be understood through the lens of physics. However, increasing attention is being paid to the
 87 interpretability of ML models (e.g., McGovern et al., 2019; Toms et al., 2020; Samek et al., 2021), and there are

88 existing methods to provide greater insight into the way physical information is propagated through them (e.g.,
89 attention maps, which identify the regions in spatial input data that have the greatest impact on the output field, and
90 ablation studies, which involve comparing reduced data sources and/or models to the original models that have full
91 access to available data, to gain insight into the models).

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

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

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

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

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

- 111 • For ML parametrization schemes:
 - 112 ○ a speedup of the representation of a subgrid-scale process vs. when run with a traditional
 - 113 parametrization scheme. This can make the difference between the scheme being cost-effective to
 - 114 run or not - when it is not cost-effective the process usually needs to be represented with a static
 - 115 forcing or boundary condition file,
 - 116 ○ a speedup of the model vs. when run with traditional parametrization schemes,
 - 117 ○ improved representation of sub-grid process(es) over traditional parameterization schemes, as
 - 118 measured by metrics appropriate to the situation,
 - 119 ○ improved overall accuracy/skill of the model when run with traditional parametrization schemes,
 - 120 ○ insight into physical processes not provided by current numerical models or theory.
- 121 • For full ML models:
 - 122 ○ a speedup of the model vs. an appropriate numerical model control,
 - 123 ○ improved overall accuracy/skill of the model vs. an appropriate numerical model control,
 - 124 ○ skillful prediction to greater lead times than an appropriate numerical model control,

125 ○ insight into physical processes not provided by current numerical models or theory

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

130 The scope of this review is deliberately limited to the application of ML within numerical weather and climate models
131 or for their replacement. This is done to keep the length of this review manageable. ML has enormous utility for other
132 aspects of the forecast value chain such as observation quality assurance, data assimilation, model output
133 postprocessing, forecast/product generation, downscaling, impact prediction, decision support tools, etc. A review of
134 the application of, and progress in, ML in these areas would be of great value but is outside the scope of this review
135 and is left to other work. Molina et al. (2023) have provided a very useful review of ML for climate variability and
136 extremes which is highly complementary to this review. They draw similar lines of delineation in the earth system
137 modelling (ESM) value chain to those mentioned above; describing them as “initializing the ESM, running the ESM,
138 and postprocessing ESM output”. They examine each of these steps in turn, with a focus on the prediction of climate
139 variability and extremes. Here we take a different approach, focusing on one part of the value chain (running the
140 ESM), but looking in more detail at this one part. Additionally, here we consider climate modelling in the context of
141 multiyear and free-running multidecadal simulations, but exclude the topic of ML for climate change projections,
142 climate scenarios, and multi-sector dynamics. This is again in the interests of ensuring the scope of the review is
143 manageable, rather than because these topics are not worthy of review. On the contrary, a review dedicated to the
144 utility of machine learning in this area would be of enormous value to the community, but cannot be adequately
145 explored here. A brief introduction to key ML architectures and concepts, including suggested foundational reading,
146 is also provided to aid readers who are unfamiliar with the subject.

147 The remainder of this review is structured as follows: In Section 2 an introduction to the two ML architectures most
148 prevalent in the review is provided, followed by a suggested methodological approach to applying ML to a problem,
149 and finally a brief overview of some of the major ML architectures and algorithms. With this background in place, the
150 application of ML in weather and climate modelling is explored in the following five sections: In Section 3, ML use
151 in sub-grid parametrization and emulation, along with tools and challenges specific to this domain, are covered.
152 Zooming out from subgrid-scale to processes resolved on the model grid, in Section 4 the application of ML for the
153 partial differential equations governing fluid flow is reviewed. Expanding scope further again to consider the entire
154 system, the use of ML for full model replacement or emulation is reviewed in Section 5. In Section 6 the growing field
155 of physics constrained ML models is introduced, and in Section 7 a number of topics tangential to the main focus of
156 this review are briefly mentioned. Setting the work covered in the previous sections in a broader context, a review of
157 the history of, and progress in, ML outside of the fields of weather and climate science is presented in Section 8. In
158 Section 9 some practical considerations for the integration of ML innovations into operational and climate models are
159 discussed, followed by a short introduction to some of the ethical considerations associated with the use of ML in
160 weather and climate modelling in Section 10. In Section 11, some future research directions are speculated on, and
161 some suggestions are made for promising areas for progression. Finally, a summary is presented in Section 12, and a

162 Glossary of Terms is provided after the final Section to aid the reader in their understanding of key concepts and
163 words.
164

165 **2. A Quick Introduction to Machine Learning**

166 While the scope of this paper is a review of ML work directly applicable to weather and climate modelling, an abridged
167 introduction to some key fundamental ML concepts is provided here to aid the reader. Suggested starting points for
168 interested readers, including guidance on the utility of different model architectures and algorithms, and the
169 connections between different applications and approaches, are as follows:

- 170 • Hsieh (2023) provides a thorough textbook on environmental data science including statistics and machine
171 learning
- 172 • Chase et al (2022a, 2022b) provide an introduction to various machine learning algorithms with worked
173 examples in a tutorial format and an excellent on-ramp to ML for weather and climate modelling
- 174 • Russell & Norvig (2021) provide a comprehensive book regarding artificial intelligence in general
- 175 • Goodfellow et al. (2016) provide a well-regarded book on deep learning theory and modern practise
- 176 • Hastie et al. (2009) provide a book on statistics and machine learning theory

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

179 The majority of ML methods which have found traction in weather and climate modelling were first developed in
180 fields such as computer vision, natural language processing and statistical modelling. Few, if any, of the methods
181 mentioned in this paper could be considered unique to weather and climate modelling, however, they have in many
182 cases been modified to a greater or lesser extent to suit the characteristics of the problem. In this review, the term
183 algorithm refers to the mathematical underpinnings of a machine learning approach. By this definition, decision trees
184 (DTs), NNs, linear regression and Fourier transforms are examples of algorithms. The two most relevant algorithms
185 for this review are DTs and NNs. Many ML algorithms can be thought of as optimizing a nonlinear regression, with
186 deep learning utilizing an extremely high-dimensional model. There is no consensus on the definition of ML, with the
187 term encompassing relatively large or small topical domains depending on who is asked. A good rule of thumb,
188 however, is that any iterative computational process that seeks to minimize a loss function or optimize an objective
189 function can be considered to be a form of ML. Some of the chief concerns in machine learning are generalizability
190 of the models, how to train (optimise the variables of) the model, and how to ensure robustness. The inputs and outputs
191 of machine learning models are the often same as physical models or model components. The term architecture in
192 machine learning refers to a specific way of utilizing an algorithm to achieve a modelling objective reliably. For
193 example, the U-Net[†] architecture is a specific way of laying out a NN which has proven effective in many applications.
194 The extreme gradient boosting decision tree[†] architecture is a specific way of utilizing DTs which has proven reliable
195 and effective for an extraordinary number of problems and situations and is an excellent choice as a first tool to
196 experiment with machine learning.

197 A major current focus of ML research in the context of weather and climate modelling is new NN-based architectures
198 and algorithms, and improved training regimes. Many other algorithms have been and continue to be employed in
199 machine learning more broadly, but are not pertinent to this review.
200 A key point for ML researchers to be aware of is the critical importance of approaching model training carefully.
201 There are many pitfalls which can result in underperformance, unexpected bias or misclassification. For instance,
202 adversarial examples[†] can occur ‘naturally’, and systems which process data can be subject to adversarial attack[†]
203 through the intentional supply of data designed to fool a trained network.

204 **2.1. Introduction to Neural Networks**

205 NNs can be regarded as universal function approximators (Hornik et al., 1989; see also Lu et al., 2019). Further, NN
206 architectures can theoretically be themselves modelled as a very wide feed-forward[†] NN with a single hidden layer.
207 A Fourier transform is another example of a function approximator, although it is not universal since not all functions
208 are periodic. NNs can therefore theoretically be candidates for accurate modelling of physical processes, although in
209 practise they cannot always reliably interpolate beyond their training envelope and as such may not generalize to new
210 regimes. ML models are typically introduced in the literature as being either classification[†] or regression[†] models, and
211 either supervised[†] or unsupervised[†].

212 The mathematical underpinning of a NN can be considered distinctly in terms of its evaluation[†] (i.e., output, or
213 prediction) step and its training update step. The prediction step can be considered as the evaluation of a many-
214 dimensional arbitrarily complex function.

215 The simplest NN is a single-input, single node network with a simple activation[†] function. A commonly used activation
216 function for a single neuron is the sigmoid function, which helpfully compresses the range between 0 and 1 while
217 allowing a nonlinear response. A classification model will employ a threshold to map the output into the target
218 categories. A regression model seeks to optimize the output result against some target value for the function. Larger
219 networks make more use of linear activations and may utilise heterogenous activation function choices at different
220 layers.

221 Complex NNs are built up from many individual nodes, which may have heterogenous activation functions and a
222 complex connectome[†]. The forward pass[†], by which inputs are fed into the network and evaluated against activation
223 functions to produce the final prediction, uses computationally efficient processes to quickly produce the result.

224 The training step for a NN is far more complex. The earliest NNs were designed by hand rather than through
225 automation. The training step applies a back-propagation[†] algorithm to apply adjustment factors to the weights[†] and
226 biases[†] of each node based on the accuracy of the overall prediction from the network.

227 Training very large networks was initially impractical. Both hardware and architecture advances have changed this,
228 resulting in the significant increase in application of NNs to practical problems. Most NN research explores how to
229 utilize different architectures to train more effective networks. There is little research going into improving the
230 prediction step as the effectiveness of a network is limited by its ability to learn rather than its ability to predict. Some

231 research into computational efficiency is relevant to the predictive step. NNs can still be technically challenging to
232 work with, and a lot of skill and knowledge are needed to approach new applications.

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

- 234 • Small, fully-connected networks, which are less commonly featured in recent publications but are still
235 effective for many tasks and are still being applied and may well be encountered in practice
- 236 • Convolutional[†] architectures, first applied to image content recognition, which match the connectome of the
237 network to the fine structure of images in hierarchical fashion to learn to recognize high-level objects in
238 images
- 239 • Recurrent token-sequence architectures, first applied to natural language processing, generation and
240 translation; applicable to any time-series problem. Now also applied to image and video applications, and
241 mixed-mode applications such as text-to-image or text-to-video
- 242 • Transformer architectures[†], based on the attention mechanism[†] to provide a non-recurrent architecture which
243 can be trained using parallelized training strategies. This allows larger models to be trained. Originally
244 developed for sequence prediction and extended to image processed through vision transformer architectures.

245 **2.2. Introduction to Decision Trees**

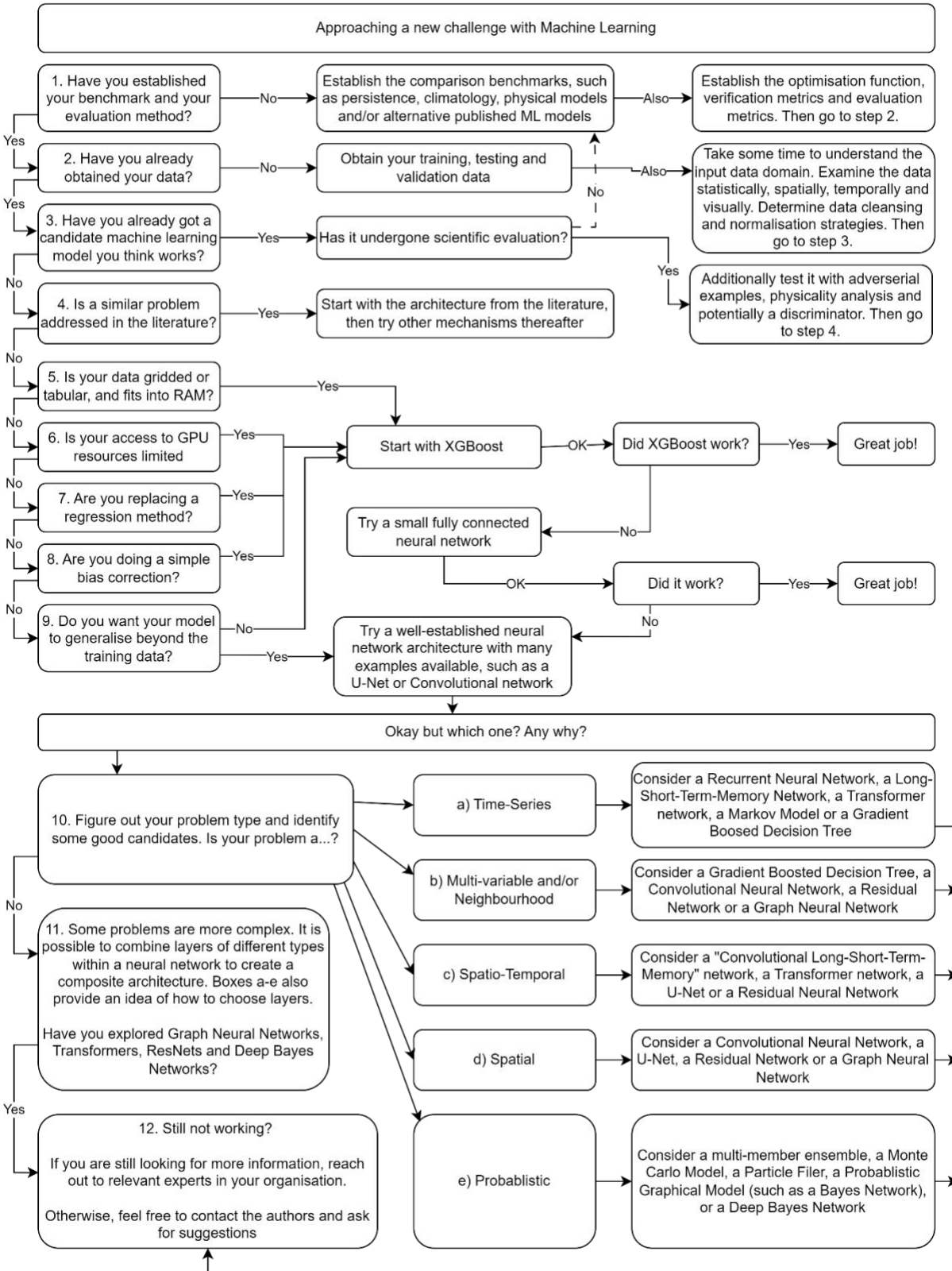
246 DTs are a series of decision points, typically represented in binary fashion based on a simple threshold. A particular
247 DT of a particular size maps the input conditions into a final 'leaf' node which represents the outcome of the decisions
248 up to that point.

249 A random forest[†] (RF) is the composition of a large number of DTs assembled according to a prescribed generation
250 scheme, which are used as an ensemble. A gradient boosted decision tree (GBDT) is built up sequentially, where each
251 subsequent decision tree attempts to model the errors of the stack of trees built up thus far. This approach outperforms
252 RFs in most cases.

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

259 **2.3. Methodologies for Machine Learning**

260 It is challenging to provide simplified advice for how to approach problem-solving in ML. There are few strict
261 theoretical reasons to choose any one of the variety of architectures which are available. The authors would also
262 caution against assuming that results in the literature are the product of a detailed comparison of alternative
263



264
265

Figure 2: A methodological flowchart illustrating a suggested approach to applying ML to a research problem.

266 architectures, or assuming that a deep learning approach is going to be easy or straightforward. It will often be the
267 case that multiple machine learning architectures may be similarly effective, and determining the optimal
268 architecture is likely to involve extensive iteration. Any specific methodology is also likely to reflect the intuitions
269 (or biases), knowledge, and background of the authors of that methodology.

270 Nonetheless, there is an appetite from many scientists for reasonable ways to 'get started' and to provide some
271 assistance for practical decision-making, particularly if approaching the utilization of machine learning for the first
272 time or in a new way. Figure 2 provides a set of suggested steps and decision points to help readers approach a new
273 challenge with ML.

274 The flowchart presented in Figure 2 provides an overview of methodological steps that can be taken when using ML
275 to solve a problem, however it does not give much insight into the pros and cons of the common ML architectures
276 available and used in the literature. Table 1 provides a brief summary of the major ML architectures and algorithms
277 used by the studies cited in this review and gives a short note on some of their pros and cons. This table is not
278 exhaustive, and readers are strongly encouraged to use it as a starting point for further exploration, rather than a
279 definitive guide. The relative strengths and weakness of each ML architecture can be subtle, and highly dependent on
280 the use case, their application, and their tuning. Establishing a good understanding of the ML architecture being used
281 is a critical step for any scientist intending to delve into ML modelling. Interested readers should also refer to Chase
282 et al (2022b), where a similar table is presented that covers a wider variety of traditional methods but fewer neural
283 network approaches.

285 An increasingly diverse array ML architectures are being applied to an ever-growing variety of challenges. These
286 architectures all have sub-variants and ancestor architectures which may not be represented, all of which may be found
287 to be of use for weather and climate modelling applications. Other concerns, such as data normalization, training
288 strategies, and capturing physicality become as relevant as the choice of architecture once a certain level of
289 performance is achieved.

290 Figure 3 shows a summary of the ML architectures and algorithms used by the studies cited in this review, including
291 the number of times each architecture is used. It can be seen from this that the two most frequently used general
292 categories of architecture are Fully Connected NNs (FCNNs) and Convolutional NNs (CNNs) of various sub-types.
293 However, some of the most significant recent research findings come from new architectures which by definition
294 cannot have wide adoption yet (these are grouped under the 'Mixed/custom NN' category in Figure 3).

295 In some cases, little justification is given for the ML architecture used in a study, and readers are therefore cautioned
296 against using the relative popularity of a particular ML architecture in the literature as a guide for its suitability for a
297 given task.

298 Furthermore, ML models increasingly use a mix of different algorithms and architectures. For example, a common
299 combination is fully-connected NN layers, convolutional NN layers, and LSTM layers. For the purposes of Figure 3,
300 the authors have endeavoured to categorise the ML architectures used in the studies in this review as accurately as
301 possible, with complex architectures being placed in the "Mixed/custom NN" category, however, where an
302 architecture was mostly but not entirely aligned with a single category, it was placed in that category. For example,
303

Approach	Description	Pros	Cons
Simple regression techniques	Includes linear regression and logistic regression. See Chase et al. (2022b) for more detail.	Explainable and well-understood.	Can only capture simple relationships.
Decision Tree	Consists of a series of branching decisions, culminating in a number of decision 'leaves'. The decision points are trainable. Provides the basis for understanding more complex decision tree and regression tree approaches.	Easily explainable. Computationally tractable and fast	Unable to fully model complex problems. Cannot make predictions outside the training envelope.
Random Forest (RF)	A random forest consists of many decision trees, which form an ensemble and the average result is taken. The construction of the trees uses randomness.	Versatile and effective. Computationally tractable and fast. Allows focus on the input variables rather than on process or model definition.	Usually performs slightly less well than gradient boosted decision trees.
Gradient Boosted Decision Trees (GBDT)	Akin to Random Forests, however each additional member is used to predict the residual error of the ensemble so far. Is often sufficient for a given problem, and should thus be considered as a baseline for measuring more complex ML models against.	A highly versatile and reliable approach. Computationally tractable and fast. Allows focus on the input variables rather than on process or model definition. Feature importance plots can guide intuition.	Has practical limitations at scale due to large memory requirements at inference time. Limited ability to simulate complex systems compared to other ML approaches such as NNs. Cannot make predictions outside the training envelope without customized leaves.

Vector Machines	Support Vector Machines (SVMs) and Relevance Vector Machines (RVMs) are supervised models used for regression and classification. RVMs have the same functional form as SVMs, but are a probabilistic classification based on Bayesian inference. Vector Machines seek to define the optimal division between classes by finding the hyperplanes which have the largest distance to the nearest training-data point of any class.	Can be used for similar problems as GBDTs. Computationally efficient and often effective. Mathematically appealing. Capable of modelling nonlinear functions.	Now less-used compared to random forests and GBDTs.
Single neuron	See Chase et al. (2022b) for a description of the structure of a perceptron. Forms the conceptual and structural basis for all NN architectures.	Unused in practice outside of a larger NN architecture.	Unable to model most problems in isolation.
Fully-Connected feed-forward Neural Network (FCNN)	Consists of multiple layers of neurons, with each neuron being connected to every neuron in the subsequent layer. Still quite widely used in weather and climate modelling, in spite of declining use in other machine learning domains. Is often sufficient and should be considered as a baseline for measuring more complex architectures against.	Effective for applications such as parametrization scheme emulation and PDE solver preconditioning. Relatively simple to work with. Computationally tractable.	Unable to effectively train beyond a certain size or depth, and thus is increasingly being replaced with more complex architectures as ML moves to deeper NNs.

Bayesian networks	A system (probabilistic graphical model) comprised of nodes which together predict both an expected value and a likelihood. Each node is associated with a probability function that provides a probability (or distribution) of the variable represented by the node.	Effective for refining an expert or knowledge-based model by incorporating additional observations. Capable of dealing with both semantic concepts and physical processes.	Determining an optimal model can be challenging and training times are prohibitive for large networks.
Deep Bayesian Networks	Deep Bayesian techniques attempt to capture the model complexity of deep neural networks while retaining the ability to predict a distribution of outcomes, a probabilistic model and a clear information-theoretical bases.	Used to obtain a more realistic expression of uncertainty. Effective in modelling where causal relationships aren't understood.	Not as well explored as neural networks in recent literature.
Convolutional Neural Network (CNN)	Involves convolving a (usually 2D image, but can also be 1D temporal, for example) input field with a filter function (often a top hat function [†]) to extract features on different spatial scales. Conceptually useful in understanding how a neural network can build up an abstract or 'big picture' definition of a process in its hidden layers by assembling fine-scale features.	The go-to network for image-based problems. Proven effective on many problems and is well-covered in the literature.	May require more significant hardware such as a modern GPU.
Residual Neural Network (ResNet)	ResNets are a form of CNN including skip connections, whereby the inputs of a number of convolutional layers are appended to the outputs of those layers to retain information lost through the weights in the convolutional layers. These skip connections make it possible to train much deeper convolutional networks than would be possible otherwise.	Allows very deep networks to be efficiently trained. Allows an iterative build-up of network size by experimenting with the number of residual layers. Could be a good choice to couple with physically interpretable layers.	Somewhat more computationally costly than other deep architectures.

<p>U-Net</p>	<p>Derives its name from the shape of the network as it is commonly shown diagrammatically (it forms a “U” shape).</p> <p>Consists of a series of downsampling convolutional layers, each of which further abstracts the information in the inputs (forming the first half of the “U”). These are then upsampled again to the original resolution of the input data (forming the second half of the “U”). Each downsampling step has its output appended to the input of the corresponding upsampling step (a form of skip connection).</p>	<p>Effective for many purposes and widely used in classification and image segmentation. Has also seen uptake for nowcasting applications and prediction of multiyear timescale ocean variables.</p>	<p>No serious drawbacks. Has somewhat given way to more complex architectures recently</p>
<p>Deep Operator Network (DeepONet)</p>	<p>A NN which is designed to learn the mappings between inputs and outputs of the mathematical operators underpinning processes, rather than directly predicting the outputs of the processes themselves. Was developed in the context of fluid dynamics and differential operators.</p> <p>An important theoretical component of the Adaptive Fourier Neural Operator used in FourCastNet (Pathak et al., 2022).</p>	<p>Provides a strong theoretical basis for learning the underlying function space of a data set.</p> <p>Highly effective for fluid dynamics and idealized systems.</p> <p>Can retain the properties of the learned operators. For example, can exhibit translational and scale invariance where that property holds for the operator in question.</p>	<p>Conceptually not straightforward.</p> <p>Requires strong mathematical and machine learning expertise to apply effectively to new challenges.</p>

<p>Graph Neural Network (GNN)</p>	<p>Models data as a set of interconnected nodes and edges (as opposed to assuming data is on a regular grid). Underpins Keisler (2022) and GraphCast (Lam et al., 2022)</p>	<p>Does not require data to be on a grid or distributed in a uniform manner. Capable of incorporating teleconnections, nonlocal relationships, and other complex variable relationships.</p>	<p>Costly to train.</p>
<p>Discriminator</p>	<p>A NN is trained to discriminate between two examples and identify the “real” one. Is used to estimate whether a sample is from the observations or the model. Forms one part of a GAN.</p>	<p>Can be used in place of a manually-defined loss function to train without over-emphasizing any individual metrics or variables. Can be used as an effective loss function when training Can be used independently to evaluate model realism. Comes closest to human subjective evaluation of image quality.</p>	<p>Is more likely to require more machine learning domain knowledge to resolve issues.</p>
<p>Generative Adversarial Network (GAN)</p>	<p>Combines a generator network with a discriminator and trains them in an adversarial manner: the discriminator tries to differentiate the generator from ground truth, the generator tries to trick the discriminator. Eventually the discriminator can't differentiate the generator from ground truth. May be part of a multi-phase training strategy in order to improve realism after initial optimization.</p>	<p>Produce results which prioritize realism over accuracy (could also be a con). Is less prone to the blurring that results from training to simpler loss functions and thus can be more effective in producing sharp images and predicting statistical extremes.</p>	<p>Increases training costs. Favors a ‘good looking’ answer over a correct answer. Can be difficult to train as the generator and discriminator must be kept balanced (one can outperform the other leading to mode collapse – a false minima).</p>

<p>Recurrent Neural Network (RNN)</p>	<p>Any neural network where the output of previous predictions are provided to a sequence-based model. Multiple sub-types of the RNN exist.</p>	<p>A simple RNN design can model many problems effectively.</p> <p>A recurrent architecture allows access to and inspection of the belief state at each iteration.</p>	<p>Recurrent approaches can accumulate errors quickly.</p> <p>Relationships which act over longer time-frames or distances than the recurrence length may not be captured.</p> <p>Choosing the length of the sequence may be a challenge.</p>
<p>Long Short Term Memory (LSTM) Network</p>	<p>Contains modified neurons with a memory component and the ability to retain or forget information. Is applied to sequence inputs and can learn the sequential scales in which information is encoded (e.g., what timescales in a timeseries are pertinent for future prediction).</p> <p>Has been combined with the ideas underpinning CNNs to create Convolutional LSTMs (ConvLSTM), which fit for both timescales of relevance and spatial features of relevance.</p>	<p>An effective alternative to a recurrent network which has proven very good at modelling time-series.</p> <p>A proven and effective mechanism for dimensionality reduction to allow the training of large networks.</p>	<p>May not include spatial relationships (unless it's a ConvLSTM), and may be more complex than needed for some problems.</p> <p>Less explainable than an attention mechanism.</p> <p>Has a bias towards closer points in a sequence (e.g., will be biased towards the recent past over a longer timescale in time series prediction).</p>

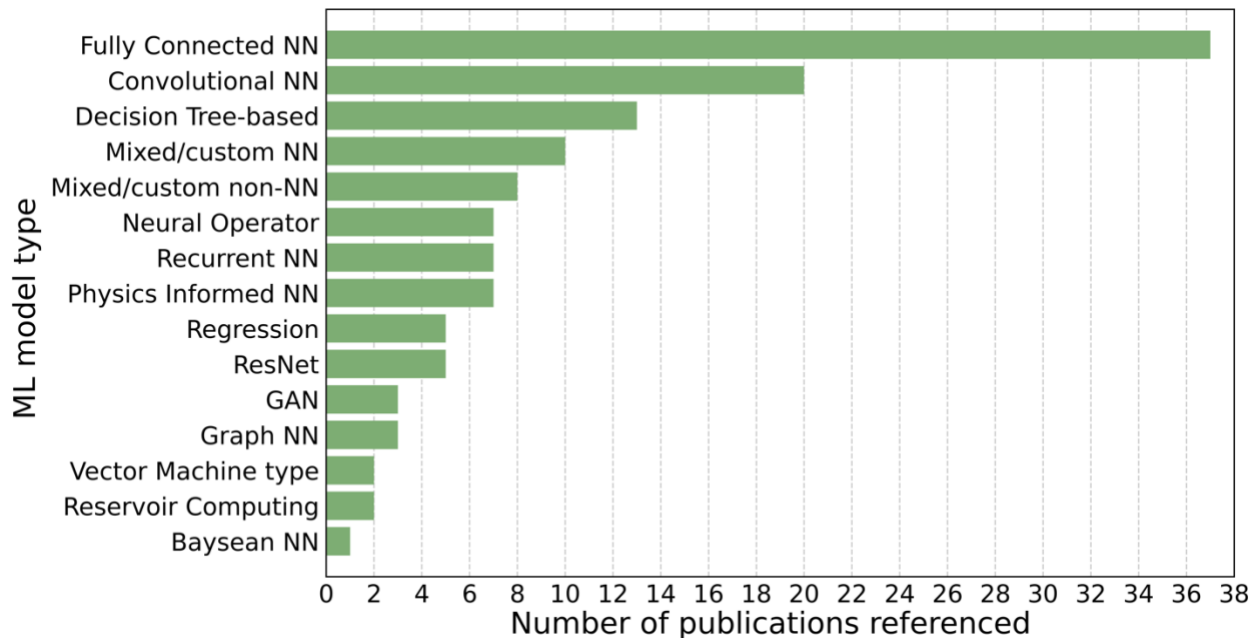
<p>Attention Mechanism</p>	<p>Often used in conjunction with other architectures as a feature extraction/dimensionality reduction method.</p> <p>A NN is trained to learn the degree of importance of each input datapoint on each other one in a sequence.</p> <p>Attention mechanism-based NNs are rapidly overtaking LSTMs as the method of choice for modelling sequence-based information.</p>	<p>Unlike LSTMS, attention mechanisms are not biased towards relationships between near points in a sequence.</p> <p>Rather, attention mechanisms treat all points in an input sequence equally and retain the learned attention mappings between each point.</p> <p>In the context of weather and climate modelling, the learned attention mappings between points can be a useful tool for assessing the degree to which a NN has learned physically realistic teleconnections.</p>	<p>More costly to train than an LSTM for the same problem because attention mechanisms have more free parameters.</p>
<p>Transformer</p>	<p>The transformer architecture combines an attention mechanism with an autoregressive approach whereby each previously predicted step in a sequence is an input into the prediction of the next step.</p> <p>Transformer architectures underpin the current generation of language models such as ChatGPT.</p> <p>Transformers are now often included as part of other architectures for input dimensionality reduction.</p>	<p>A proven and effective mechanism for dimensionality reduction to allow the training of large networks.</p> <p>While the uptake of transformer architectures in weather and climate modelling is still small, their impressive performance for sequence prediction suggests they could have great for the field.</p>	<p>Transformers can be difficult to train due to a tendency to overemphasize the recurrent component of the network over new inputs in the early stages of training.</p>

304
305
306
307
308
309
310
311

Table 1: A summary of major ML architectures and algorithms used by the studies cited in this review. Interested readers should also refer Chase et al (2022b) where a similar table is presented that covers a wider variety of traditional methods but fewer neural network approaches.

an LSTM model with a small number of feed-forward layers would be categorised as a Recurrent NN. Since many contemporary ML models combine multiple architectural elements and algorithms into the one model, it is somewhat of an oversimplification to consider each of these in isolation, and while starting with a simple model design with a

312 limited selection of layer types is advisable to aid interpretability, there is no reason they cannot be combined or used
 313 in conjunction with each other if this improves the performance of the model.
 314 Adapting, optimizing and debugging issues with machine learning systems can be very complex (especially so for
 315 large NNs), and is likely to require both machine learning expertise and domain knowledge (i.e. scientific knowledge).
 316 XGBoost provides the ability to generate chart showing the importance of the features in the model which can be very
 317 helpful. Shapley Additive Explanations (Lundberg and Lee 2017) can provide insights into feature importance for any
 318 model including NNs.
 319



320
 321 Figure 3: A count of the ML architectures and algorithms used by the studies cited in this review. As with Figure 1, this figure
 322 includes all references from this review except for: seminal ML papers that are on new ML methods (e.g., foundational ML papers),
 323 review papers, any paper cited that concerns a topic which is out of scope (e.g., nowcasting), and any other paper which does not
 324 present a new method directly applicable to weather and climate modelling. The full table of citations is provided in the appendix.
 325

326 3. Sub-grid parametrization and emulation

327 Subgrid-scale processes in numerical weather and climate models are typically represented via a statistical
 328 parameterization of what the macroscopic impacts of the process would be on resolved processes and parameters.
 329 These are commonly referred to as parameterization schemes, and can be very complex and relatively computationally
 330 costly. For example, in the European Centre for Medium-Range Weather Forecast’s (ECMWF) Integrated Forecasting
 331 System (IFS) model they account for about a third of the total computational cost of running the model (Chantry et al.
 332 2021b). They also require some understanding of the underlying unresolved physical processes. Examples of subgrid-
 333 scale processes which are typically currently parameterized in operational systems include gravity wave drag,
 334 convection, radiation, subgrid-scale turbulence, and cloud microphysics. As additional complexity (for example

335 representation of aerosols, atmospheric chemistry, land surface processes, etc.) is added to numerical models, the
336 computational cost will only increase.

337 ML presents an alternative approach to representing subgrid-scale processes, either by emulating the behavior of an
338 existing parametrization scheme, emulating the behavior of sub-components of the scheme, by replacing the current
339 scheme or sub-component entirely with an ML-based scheme, or by replacing the aggregate effects of multiple
340 parametrization schemes with a single ML model.

341 ML emulation of existing schemes or sub-components has the advantage of maintaining the status quo within the
342 model; no or minimal re-tuning of the model should be required since the ML emulation is trained to replicate the
343 results of an already-tuned-for scheme. Because of this, the main benefit of this approach is that it reduces the
344 computational cost of running the parametrization scheme. On the other hand, full replacement of an existing
345 parameterization scheme or sub-component with an ML alternative has the potential to be both computationally
346 cheaper and also an improvement over the preceding scheme.

347 In the following subsections, a review of the literature on aspects of ML for the parametrization and emulation of
348 subgrid-scale processes is presented.

349 **3.1. Early work on ML parametrization and ML emulations**

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

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

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

367 **3.2. ML for coarse graining**

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

371 This has proven to be a popular method for developing ML-based parametrization schemes. Brenowitz & Bretherton
372 (2018) used a near-global aqua planet simulation run at 4 km grid length to train a NN to represent the apparent sources
373 of heat and moisture averaged onto 160 km² grid boxes. They then tested this scheme in a prognostic single column
374 model and showed that it performed better than a traditional model in matching the behavior of the aqua planet
375 simulation it was trained on. Brenowitz & Bretherton (2019) built on this work by training their NN on the same global
376 aqua-planet 4 km simulation, but then embedded this scheme within a coarser resolution (160 km²) global aqua planet
377 GCM. Embedding NNs within GCMs is challenging because feedbacks between NN and GCM components can cause
378 spatially extended simulations to become dynamically unstable within a few model days. This is due to the inherently
379 chaotic nature of the atmosphere in the GCM responding to inputs from the NN which cause rapidly escalating
380 dynamical instabilities and/or violate physical conservation laws. The authors overcame this by identifying and
381 removing inputs into the NN which were contributing to feedbacks between the NN and GCM (Brenowitz et al. 2020),
382 and by including multiple time steps in the NN training cost function. This resulted in stable simulations which
383 predicted the future state more accurately than the course resolution GCM without any parametrization of subgrid-
384 scale variability, however the authors do observe that the mean state of their NN-coupled GCM would drift, making
385 it unsuitable for prognostic climate simulations.

386 Rasp et al. (2018) trained a deep NN[†] to represent all atmospheric subgrid processes in an aquaplanet climate model
387 by learning from a multiscale model in which convection was treated explicitly. They then replaced all sub-grid
388 parameterizations in an aquaplanet GCM with the deep NN, and allowed it to freely interact with the resolved
389 dynamics and the surface-flux scheme. They showed that the resulting system was stable and able to closely reproduce
390 not only the mean climate of the cloud-resolving simulation but also key aspects of variability in prognostic multiyear
391 simulations. The authors noted that their decision to use deep NNs was a deliberate one, because they proved more
392 stable in their prognostic simulations than shallower NNs, and they also observed that larger networks achieved lower
393 training losses. However, while Rasp et al. (2018) were able to engineer a stable model that produced results close to
394 the reference GCM, small changes in the training dataset or input and output vectors quickly led to the NN producing
395 increasingly unrealistic outputs and causing model blow-ups (Rasp 2020). Consistent with this, Brenowitz &
396 Bretherton (2019) report that they were unable to achieve the same improvements in stability with increasing network
397 layers found by Rasp et al. (2018).

398 **3.3. Overcoming instability in ML emulations and parametrizations**

399 O’Gorman & Dwyer (2018) tackled the instabilities observed in NN-based approaches to subgrid-scale
400 parameterization by employing an alternative ML method; Random Forests (RFs; Breiman 2001; Tibshirani &
401 Friedman 2001). The authors trained a RF to emulate the outputs of a conventional moist convection parametrization

402 scheme. They then replaced the conventional parameterization scheme with this emulation within a global climate
403 model, and showed that it ran stably and was able to accurately produce climate statistics such as precipitation
404 extremes without needing to be specially trained on extreme scenarios. RFs consist of an ensemble of DTs, with the
405 predictions of the RF being the average of the predictions of the DTs which in turn exist within the domain of the
406 training data. RFs thus have the property that their predictions cannot go outside of the domain for their training data,
407 which in the case of O’Gorman & Dwyer (2018) ensured conservation of energy and nonnegativity of surface
408 precipitation (both critically important features of the moist convection parametrization scheme) were automatically
409 achieved. A disadvantage of this method however is that it requires considerable memory when the climate model is
410 being run to store the tree structures and predicted values which make up the RF.

411 Yuval & O’Gorman (2020) extended on the ideas in O’Gorman & Dwyer (2018), switching from emulation of a single
412 parametrization scheme to emulation of all atmospheric sub grid processes. They trained an RF on a high-resolution
413 three-dimensional model of a quasi-global atmosphere to produce outputs for a course-grained version of the model,
414 and showed that at course resolution the RF can be used to reproduce the climate of the high-resolution simulation,
415 running stably for 1000 days.

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

428 **3.4. From aquaplanets to realistic land-ocean simulations**

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

433 Mooers et al. (2021) represents a subsequent example of an ML emulation of atmospheric fields with realistic
434 geographical boundary conditions, where the authors developed feed-forward NNs to super-parametrize subgrid-scale
435 atmospheric parameters and forced a realistic land surface model with them. Super-parametrization is distinct from
436 traditional parameterization in that it relies on solving (usually simplified) governing equations for subgrid-scale

437 processes rather than heuristic approximations of these processes. They employed automated hyperparameter
438 optimization[†] to investigate a range of neural network architectures across ~250 trials, and investigated the statistical
439 characteristics of their emulations. While the authors found that their NNs had a less good fit in the tropical marine
440 boundary layer, attributable to the NN struggling to emulate fast stochastic signals in convection, they also reported
441 good skill for signals on diurnal to synoptic timescales.

442 Brenowitz et al. (2022) sought to address the challenge of emulating fast processes. They used FV3GFS (Zhou et al.,
443 2019; Harris et al., 2021; a compressible atmospheric model used for operational weather forecasts by the US National
444 Weather Service) with a simple cloud microphysics scheme included to generate training data and used this to train a
445 selection of ML models to emulate cloud microphysics processes, including fast phase changes. They emulated
446 different aspects of the microphysics with separate ML models chosen to be suitable to each task. For example, simple
447 parameters were trained with single-layer NNs, while parameters which are more complex spatially were trained with
448 RNNs (e.g., rain falls downwards and not upwards, so it is sequential in timesteps through the atmosphere – a feature
449 which can be represented by an RNN). They then embedded their ML emulation in FV3GFS. They found that their
450 combined ML simulation performed skillfully according to their chosen metrics, but had excessive cloud over the
451 Antarctic Plateau.

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

453 **3.5. Testing with prognostic long-term simulations**

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

461 Chantry et al. (2021b) trained a NN to emulate the non-orographic gravity wave drag parameterization in the ECMWF
462 IFS model (specifically cycle 45R1, ECMWF, 2018) and were able to run stable, accurate simulations out to 1 year
463 with this emulation coupled to the IFS. While the authors note that RFs have been shown to be more stable (e.g.,
464 O’Gorman & Dwyer (2018) and Yuval & O’Gorman (2020), as described above, and Brenowitz et al. (2020)), they
465 chose to focus on NNs since they have lower memory requirements and therefore promise better theoretical
466 performance. The authors assessed the performance of their emulation in a realistic GCM by coupling the NN with
467 the IFS, replacing the existing non-orographic gravity wave drag scheme, and performed 120 hour, 10 day, and 1 year
468 forecasts at ~25 km resolution in a variety of model configurations. The authors showed that their emulation was able
469 to run stably when coupled to the IFS for seasonal timescales, including being able to reproduce the descent of the
470 Quasi-biennial Oscillation (QBO). Interestingly, while the authors initially aimed to ensure momentum conservation
471 in a manner similar to Beucler et al. (2021), they found that this constraint led to model instabilities and that a better

472 result was achieved without it. One possible explanation for this is that Beucler et al. (2021) assessed their NNs in an
473 aquaplanet setting. Nonetheless, Chantry et al. (2021b) noted that since their method was not identical to Beucler et
474 al. (2021), improved stability could potentially be achieved by following their method more precisely. The
475 computational cost of the NN emulation developed by Chantry et al. (2021b) was found to be similar that of the
476 existing parametrization scheme when run on CPUs, but was faster by a factor of 10 when run on GPUs due to the
477 reduction in data transmission bottlenecks.

478 The first study to successfully run stable long-term climate simulations with ML parametrizations was Wang et al.
479 (2022a), who extended on the work of Han et al. (2020) by constructing a ReNet to emulate moist physics processes.
480 They used the residual connections from Han et al. (2020) to construct NNs with good nonlinear fitting ability, and
481 filtered out unstable NN parametrizations using a trial-and-error analysis, resulting in the best ResNet set in terms of
482 accuracy and long-term stability. They implemented this scheme in a GCM with realistic geographical boundary
483 conditions and were able to maintain stable simulations for over 10 years in an Atmospheric Model Intercomparison
484 Project (AMIP)-style configuration. This was more akin to a hybrid ML-physics based model than a traditional GCM
485 with ML-based parametrization, because rather than embedding the ResNet in the model code, the authors used a NN-
486 GCM coupling platform through which the NNs and GCMs could interact through data transmission. This is in
487 contrast to the approach employed in the Physical-model Integration with Machine Learning⁴ (PIML) project and
488 Infero⁵, which are both described in Section 3.1.1. One advantage to this approach noted by the authors is that it allows
489 for a high degree of flexibility in the application of the ML component, however is likely to be less efficient than a
490 fully-embedded ML model, due to the potential for data transmission bottlenecks.

491 **3.6. Training with observational data**

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

501 **3.7. ML for super parameterization**

502 Revisiting the topic of super parametrized subgrid-scale processes introduced above, the use of ML for this approach
503 was investigated in depth by Chattopadhyay et al. (2020). The authors introduced a framework for NN-based super

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

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

504 parametrization, and compared the performance of this method against NN-based traditional parametrization (i.e.,
505 based on heuristic approximations of subgrid-scale processes) and direct super parameterization (i.e., explicitly
506 solving for the subgrid-scale processes) in a chaotic Lorenz '96 (Lorenz 1996) system that had three sets of variables,
507 each of a different scale. They found that their NN-based super parameterization outperformed direct super
508 parameterization in terms of computational cost, and was more accurate than NN-based traditional parametrization.
509 The NN-based super parameterization showed comparable accuracy to direct super parameterization in reproducing
510 long-term climate statistics, but was not always comparable for short-term forecasting.

511 **3.8. Stochastic parametrization schemes**

512 A more recent approach to the representation of subgrid-scale processes is via stochastic parameterization schemes,
513 which can represent uncertainty within the scheme. There has been less focus on replacing these schemes with ML
514 alternatives than non-stochastic schemes, however some progress has been made. Krasnopolsky et al. (2013) used an
515 ensemble of NNs to learn a stochastic convection parametrization from data from a high-resolution cloud resolving
516 model. In this case, the stochastic nature of the parametrization was captured by the ensemble of NNs. Gagne et al
517 (2020b) took a different approach, investigating the utility of generative adversarial networks (GANs) for stochastic
518 parametrization schemes in Lorenz '96 (Lorenz 1996) models. In this case, the GAN learned to emulate the noise of
519 the scheme directly, rather than implicitly representing it with an ensemble. They described the effects of different
520 methods to characterize input noise for the GAN, and the performance of the model at both weather and climate
521 timescales. The authors found that the properties of the noise influenced the efficacy of training. Too much noise
522 resulted impaired model convergence and too little noise resulted in instabilities within the trained networks.

523 **3.9. ML parametrization and emulation for land, ocean, and sea ice models**

524 Models of the atmosphere make up one component of the Earth system, however for timescales beyond a few days,
525 simulating other components of the Earth system becomes increasingly important to maintain accuracy. The
526 components which are most often included in coupled Earth system models in addition to the atmosphere are the
527 ocean, sea ice, and the land surface. Reflective of this, ML approaches to parameterization of subgrid-scale processes
528 are not limited to the atmosphere, and progress has been made in the use of ML for land, ocean and sea ice models as
529 well.

530 On the ocean modelling front, Krasnopolsky et al. (2002) presented an early application of NN for the approximation
531 of seawater density, the inversion of the seawater equation of state, and a NN approximation of the nonlinear wave-
532 wave interaction. More recently, Bolton & Zanna (2019) investigated the utility of Convolutional Neural Networks
533 (CNNs) for parametrizing unresolved turbulent ocean processes and subsurface flow fields. Zanna & Bolton (2020)
534 then investigated both Relevance Vector Machines[†] (RVMs) and CNNs for parameterizing mesoscale ocean eddies.
535 They demonstrated that because RVMs are interpretable, they can be used to reveal closed-form equations for eddy
536 parameterizations with embedded conservation laws. The authors tested the RVM and CNN parameterizations in an
537 idealized ocean model and found that both improved the statistics of the coarse resolution simulation. While the CNN

538 was found to be more stable than the RVM, the advantage of the RVM was the greater interpretability of its outputs.
539 Finally, Ross et al. (2023) developed a framework for benchmarking ML based parametrization schemes for subgrid-
540 scale ocean processes. They used CNNs, symbolic regression, and genetic programming methods to emulate a variety
541 of subgrid-scale forcings including measures of potential vorticity and velocity, and developed a standard set of
542 metrics to evaluate these emulations. They found that their CNNs were stable and performed well when implemented
543 online, but generalized poorly to new regimes.

544 Focusing instead on sea ice, Chi & Kim (2017) assessed the ability of two NN models; a fully connected NN and an
545 LSTM, to predict Antarctic sea ice concentration up to a year in advance. Their ML models outperformed an
546 autoregressive model comparator, and were in good agreement with observed sea ice extent. Andersson et al. (2021)
547 improved upon this work with their model IceNet, A U-Net ensemble model which produced probabilistic Arctic sea
548 ice concentration predictions to a 6-month lead time. The authors compared IceNet to the SEAS5 dynamical sea ice
549 model (Johnson et al., 2019) and showed an improvement in the accuracy of a binary classification of ice/no ice for
550 all lead months except the first month. Horvat & Roach (2022) used ML to emulate a parameterization of wave-
551 induced sea ice floe fracture they had developed previously, in order to reduce the computational cost of the scheme.
552 When embedded in a climate simulation, their ML scheme resulted in an overall categorical accuracy (accounting for
553 the fact that it was only called where needed) of 96.5%. However, the authors did note that since their ML scheme
554 was trained on present day sea ice conditions, it may have reduced success under different climate scenarios, and they
555 recommend retraining using climate model sea-ice conditions to account for this. Rosier et al. (2023) developed
556 MELTNET, a ML emulation of the ocean induced ice shelf melt rates in the NEMO ocean model (Gurvan et al.,
557 2019). MELTNET consisted of a melt rate segmentation task, followed by a denoising autoencoder network which
558 converted the discrete labelled melt rates to a continuous melt rate. The authors demonstrated that MELTNET
559 generalized well to ice shelf geometries outside the training set, and outperformed two intermediate-complexity melt
560 rate parameterizations, even when parameters in those models were tuned to minimize any misfit for the geometries
561 used. Given the computational cost of sea ice parametrizations is relatively high for the timescales on which sea ice
562 evolution is important (namely, seasonal to climate timescales), and given the promising results in emulating these
563 parametrizations demonstrated in the literature, ML based emulation of these schemes is a strong candidate for
564 inclusion into future dynamical coupled modelling systems.

565 Finally, considering Earth's surface, most of the focus of ML innovations in this context has focused on land use
566 classification (e.g, Carranza-García et al, 2019; Digra et al., 2022) and crop modelling (e.g., Virmodkar et al., 2020;
567 Zhang et al., 2023). The rate of publication of ML applications for land surface models has been slower, however
568 there has nonetheless been steady progress in this space in recent years. Pal & Sharma (2021) presented a review of
569 the use of ML in land surface modelling which provides an excellent primer of the state of the field to that point. They
570 include in their review an overview of land surface modelling components and processes, before reviewing the
571 literature on the use of ML to represent them. They separate their review into attempts to predict and parametrize
572 different variables or aspects of the model, including evapotranspiration (Alemohammad et al., 2017; Zhao et al.,
573 2019; Pan et al., 2020), soil moisture (Pelissier et al., 2020), momentum and heat fluxes (Leufen & Schädler, 2019),
574 and parameter estimation and uncertainty (Chaney et al., 2016; Sawada, 2020; Dagon et al., 2020). They also provide

575 a useful summary of the ML architectures that have been used in publications they have discussed. More recently, He
576 et al. (2022) developed a hybrid approach to modelling aspects of the land surface, where a traditional land surface
577 model was used to optimize selected vegetation characteristics, while a coupled ML model simulated a corresponding
578 three-layer soil moisture field. The estimated evapotranspiration from this hybrid model was compared to observations
579 and it was found that it performed well in vegetated areas but underestimated the evapotranspiration in extreme arid
580 deserts. The ready application of ML to aspects of land surface modelling, and the relative sparsity of publications in
581 this space suggests that it is a fertile domain for further research and development.

582 **3.10. ML for representing or correcting a sub-component of a parametrization scheme**

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

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

589 Meyer et al. (2022) trained a NN to account for the differences between 1D cloud effects in the European Centre for
590 Medium Range Weather Forecasting (ECMWF) 1D radiation scheme ecRad and 3D cloud effects in the ECMWF
591 SPARTACUS (SPeedy Algorithm for Radiative TrAnSfer through CloUd Sides) solver. The 1D cloud effects solver
592 within ecRad, Tripleclouds, is favored over the 3D SPARTACUS solver because it is five times less computationally
593 expensive. The authors show that their NN can account for differences between the two schemes with typical errors
594 between 20% and 30% of the 3D signal, resulting in an improvement in Tripleclouds' accuracy with an increase in
595 runtime of approximately 1%. By accounting for the differences between SPARTACUS and Tripleclouds rather than
596 emulating all of SPARTACUS, the authors were able to keep Tripleclouds unchanged within ecRad for cloud-free
597 areas of the atmosphere, and utilize the NN 3D correction elsewhere.

598 **3.11. Bridging the gap between popular languages for ML and large numerical models**

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

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

607 Wang et al. (2022a; mentioned already above) opted for method 1, developing what could be considered a hybrid ML-
608 physics based model rather than a traditional GCM with ML-based parametrization. In their study, the authors used a

609 NN-GCM coupling platform through which the NNs and GCMs could interact through data transmission. One
610 advantage to this approach noted by the authors is that it allows for a high degree of flexibility in the application of
611 the ML component, however, is likely to be less efficient than a fully-embedded ML model, due to the potential for
612 data transmission bottlenecks. This framework was then formalized by Zhong et al. (2023).
613 There are many examples where method 2 was used, such as Rasp et al. (2018), Brenowitz & Bretherton (2018),
614 Gagne et al. (2019) and Gagne et al. (2020a). The obvious disadvantage of this approach is that every change to the
615 ML model being used requires reimplementing in the Fortran, and if the aim is to test a suite of ML models, this
616 approach becomes untenable. Furthermore, this approach poses greater technical barriers for scientists developing
617 ML-based solutions for numerical model challenges, since they must be sufficiently proficient in Fortran to
618 reimplement models in it, rather than using existing user-friendly Python toolkits.
619 A solution lying somewhere between methods 2 and 3 was developed by Ott et al. (2020), who developed a Fortran-
620 Keras Bridge (FKB) library that facilitated the implementation of Keras-like[†] NN modules in Fortran, providing a
621 more modular means to build NNs in Fortran code. This however did not fully overcome the drawbacks posed by
622 method 2 on its own; implementation of layers in the Fortran is still necessary, and any innovations in the Python
623 modules being used would need to be mirrored in the Fortran library.
624 Finally, method 3 is being tackled by the Met Office in the PIML⁶ project, and by ECMWF with an application called
625 Infero⁷. These projects both seek to develop a framework which can be used by researchers to develop ML solutions
626 to modelling problems in Python, and then integrate them directly into the existing codebase of the physical model
627 (e.g., the Unified model at the UK Met Office). The approach used is to directly expose the compiled code
628 underpinning the Python modules within the physical model code.

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

630 The representation and solving of the partial differential equations (PDEs) governing the fluid flow and dynamical
631 processes in the oceans and atmosphere can be considered the backbone of weather and climate models. The solvers
632 used to find solutions to these equations are typically iterative, and must solve the dynamics-governing equations of
633 their model on every timestep and at every grid point. There has been growing interest in using ML to facilitate
634 speedups and computational cost reductions in the preconditioning and execution of these solvers. Preconditioners are
635 used to reduce the number of iterations required for a solver to converge on a solution, and usually do so by inverting
636 parts of the linear problem. Many earlier studies focused on using ML to select the best preconditioner and/or PDE
637 solver from a set of possible choices (e.g. Holloway & Chen, 2007; Kuefler & Chen, 2008; George et al., 2008; Pairs
638 & Chen, 2011; Huang et al., 2016; and Yamada et al., 2018). Ackmann et al. (2020) approached the preconditioner
639 part of the system more directly, using a variety of ML methods to directly predict the pre-condition of a linear solver,
640 rather than using a standard preconditioner. Rizzuti et al. (2019) focused on the solver, using ML to apply corrections

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

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

641 to a traditional iterative solver for the Helmholtz equation. Going a step further, a number of studies have used ML to
642 replace the linear solver entirely (Ladický et al., 2015; Yang et al., 2016; Tompson et al., 2017).

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

650 Building on this, Kochkov et al. (2021) developed a ML-based method to accurately calculate the time evolution of
651 solutions to nonlinear PDEs which used grids an order of magnitude coarser than is traditionally required to achieve
652 the same degree of accuracy. They used convolutional NNs to discover discretized versions of the equations (as in
653 Bar-Sinai et al., 2019), and applied this method selectively to the components of traditional solvers most affected by
654 coarse resolution, with each NN being equation specific. They utilized the property that the dynamics of the PDEs
655 were localized, combined with the convolutional layers of their NN enforcing translation invariance[†], to perform their
656 training simulations on small but high-resolution domains, making the training set affordable to produce. An
657 interesting feature of their training approach, which is growing in popularity, was the inclusion of the numerical solver
658 in the training loss function: the loss function was defined as the cumulative pointwise error between the predicted
659 and ground truth values over the training period. In this way, the NN model could see its own outputs as inputs,
660 ensuring an internally-consistent training process. This had the effect of improving the predictive performance of the
661 model over longer timescales, in terms of both accuracy and stability. Finally, the authors demonstrated that their
662 models produced generalizable properties (i.e., although the models were trained on small domains, they produced
663 accurate simulations over larger domains with different forcing and Reynolds number). They showed that this
664 generalization property arose from consistent physical constraints being enforced by their chosen method.

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

674 The application of ML to the solving of PDEs and the preconditioning of PDE solvers has been a fruitful avenue of
675 research to date. It has led to innovations which have proven useful even outside of the immediate field (e.g., Pathak
676 et al. 2022 adapted innovations from DeepONets to use in fully ML-based weather models - this is discussed further
677 in the next Section). This is likely in part because there are many areas of engineering and science which are active in

678 progressing relevant research, leading to a greater overall pace of innovation. ML-based PDE solvers and
679 preconditioners have not yet been tested in a physical weather and climate model. There are few theoretical reasons
680 this could not occur and, if effective, result in significant computational efficiencies for traditional physical model
681 architectures. This poses an interesting avenue for further research.

682 **5. Numerical model replacement/emulation**

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

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

688 **5.1. Early work – 1D deterministic models**

689 Work on the use of ML to predict chaotic time-domain systems initially focused on 1-D problems, including 1-D
690 Lorenz systems (e.g. Karunasinghe & Liang, 2006; Vlachas et al., 2018). Of particular interest is Vlachas et al. (2018),
691 who used Long Short-Term Memory Networks (LSTMs[†]), which are well-suited to complex time domain problems.
692 Convolutional LSTMs (ConvLSTMs), which combine convolutional layers with an LSTM mechanism, were
693 introduced in the meteorological domain by Shi et al. (2015) for precipitation nowcasting. They have since seen wide
694 adoption in other areas (e.g., Yuan et al., 2018; Moishin et al., 2021; Kelotra & Pandey, 2020). Their success in other
695 domains suggests that revisiting their utility for weather and climate modelling could be worthwhile.

696 **5.2. Moving to spatially extended deterministic ML-based models**

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

710 Weyn et al. (2019) aimed to predict a limited number of variables, focusing on the NWP to medium range time domain.
711 They trained a CNN to predict 500 hPa geopotential height and 300 to 700 hPa geopotential thickness over the
712 Northern Hemisphere to up to 14-days lead time, showing better skill out to 3 days than persistence, climatology, and
713 a dynamics-based barotropic vorticity model, but not better than an operational full-physics weather prediction model.
714 Weyn et al. (2020) then improved on this significantly, with a Deep U-Net style CNN trained to predict four variables
715 (geopotential height at 500 and 1000 hPa, 300 to 700 hPa geopotential thickness, and 2 m temperature) globally to 14
716 days lead time. A major innovation in this study was their use of a cubed-sphere grid, which minimized distortions
717 for planar convolution algorithms while also providing closed boundary conditions for the edges of the cube faces.
718 Additionally, they extended their previous work to include sequence prediction techniques, making skillful predictions
719 possible to longer lead times. Their improved model outperformed persistence and a coarse resolution comparator (a
720 T42 spectral resolution version of the ECMWF IFS model, with 62 vertical levels and ~2.8 degree horizontal
721 resolution) to the full 14 days lead time, but was not as skillful as a higher resolution comparator (a T63 spectral
722 resolution version of the IFS model with 137 vertical levels and ~1.9 degree horizontal resolution) or the operational
723 subseasonal-to-seasonal (S2S) version of the ECMWF IFS.

724 Clare et al. (2021) tackled a short falling of many of the ML weather and climate models developed to this point,
725 namely that most were deterministic, limiting their potential utility. To address this, they trained a NN to predict full
726 probability density functions of geopotential height at 500 hPa and temperature at 850 hPa at 3 and 5 days lead time,
727 producing a probabilistic forecast which was comparable in accuracy to Weyn et al. (2020).

728 Choosing to focus on improved skill rather than the question of probabilistic vs deterministic models, Rasp & Thuerey
729 (2021) developed a ResNet model trained to predict geopotential height, temperature and precipitation to 5 days lead
730 time and assessed it against the same set of physical models as Weyn et al. (2020). Their model was close to as skillful
731 as the T63 spectral resolution version of the IFS model, and had better skill to the 5 day lead time than Weyn et al.
732 (2020).

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

742 Lam et al. (2022) also used GNNs to build their ML-based weather and climate model, GraphCast. This model was
743 the most skillful ML-based weather and climate model at the time of writing this review. While the first ML-based
744 weather and climate model to claim to exceed the skill of a numerical model was Pangu-Weather (Bi et al., 2022;
745 described in greater detail in the following subsection), GraphCast exceeded the skill of both the ECMWF
746 deterministic operational forecasting system, HRES, and also Pangu-Weather. Furthermore, Lam et al. (2022) paid

747 particular attention to evaluating their model and HRES against appropriate measures, and included existing model
748 assessment scorecards from ECMWF to evaluate them. GraphCast capitalized on the ability of GNNs to model
749 arbitrary sparse interactions by adopting a high-resolution multi-scale mesh representation of the input and output
750 parameters. It was trained on the ECMWF ERA5 reanalysis archive to produce predictions of five surface variables
751 and six atmospheric variables, each at 37 vertical pressure levels, on a 0.25° grid. It made predictions on a 6-hourly
752 timestep and was run autoregressively to produce predictions to a 10-day lead time. The authors demonstrated that
753 GraphCast was more accurate than HRES on 90.0% of the 2760 variable and lead time combinations they evaluated.

754 **5.3. Ensemble generation with ML-based models**

755 A common criticism of ML approaches to weather and climate prediction is the difficulty of representing uncertainty,
756 and/or the tails of the distribution of predicted parameters. One common method to represent the range of possible
757 outcomes (including extremes) under different sources of uncertainty is through a well-calibrated ensemble of
758 predictions. There are a growing number of examples where ensemble generation is considered, many of which fall
759 into the category of full-model replacement.

760
761 Weyn et al. (2021) explored probabilistic ML predictions using an ensemble of NNs similar to the single-member NN
762 described in Weyn et al. (2020). The authors expanded the number of variables predicted from 4 to 6, and produced
763 forecasts to 6 weeks lead time - considerably longer than any comparable work at the time of writing this review. They
764 considered a variety of initial condition perturbation strategies, and explored the impact of model error by varying the
765 initial values of the model weights during training to create a multi-model ensemble. They used a combination of the
766 multi-model ensemble generation approach and initial condition perturbations to generate a ‘grand ensemble’ of 320
767 members. They used established metrics for ensemble performance such as RMSE-spread plots, and found that the
768 320-member grand ensemble combining the multi-model ensemble with initial condition perturbations performed only
769 slightly better than the multi-model ensemble alone at 14 day lead times. The skill of the ensemble mean of the system,
770 a control member, and the full ensemble were assessed against the same metrics from the ECMWF sub-seasonal to
771 seasonal (S2S) prediction system. Their grand ensemble had lower skill than the S2S system at shorter lead times, but
772 was comparable in skill at longer lead times. Their skill assessment used standard probabilistic skill measures such as
773 continuous ranked probability score and the ranked probability skill score, which are not present in the other studies
774 discussed in this Section. The next major ML model to be tested in an ensemble mode was FourCastNet, presented by
775 Pathak et al. (2022), who leveraged the work on DeepONets described in Section 4. In particular, the authors used a
776 type of DeepONet called a Fourier Neural Operator (FNO). FourCastNet produced predictions of 20 variables
777 (including challenging-to-predict variables such as surface winds and precipitation) on five vertical levels with 0.25
778 degree horizontal resolution, and had competitive skill against the ECMWF IFS to 1 week lead time. The high
779 horizontal resolution of their model enabled it to resolve extreme events such as tropical cyclones and atmospheric
780 rivers, and the speed of the model facilitated the generation of large ensembles (up to 1,000’s of members).

781 The authors explored the potential of their ensemble forecasts by generating a 100-member ensemble from initial
782 conditions perturbed with Gaussian random noise. They showed that the FourCastNet ensemble mean had lower
783 RMSE and a higher anomaly correlation coefficient than a single-value prediction at longer lead times (beyond ~3-4
784 days), although the ensemble mean performed slightly worse than the single value forecast at shorter lead times. The
785 authors attributed this relative decrease in performance at shorter lead times to the ensemble mean smoothing out fine-
786 scale features. Unfortunately, the authors did not examine the spread of the ensemble with lead time or evaluate the
787 model using probabilistic skill metrics (in contrast to Weyn et al., 2021), and while they did consider the capacity of
788 FourCastNet to predict extremes, they did not do so in an ensemble context.

789 Hu et al. (2023) improved on the relatively simple ensemble perturbation approach employed by Pathak et al. (2022)
790 in their model, a Swin (sliding window) Transformer-based Variational Recurrent Neural Network (SwinVRNN).
791 This model combined a Swin Transformer Recurrent Neural Network (SwinRNN) predictor with a Variational Auto-
792 Encoder perturbation module. The perturbation module learned the multivariate Gaussian distributions of a time-
793 variant stochastic latent variable from the training data. The SwinRNN predictor was deterministic, but could be used
794 to generate ensemble predictions by perturbing model features using noise sampled from the distribution learned by
795 the perturbation module. Unlike the approach used by Pathak et al. (2022), this strategy ensured that the perturbations
796 applied at each spatial location in ensemble generation were appropriate for the location and variable in question.
797 Furthermore, the training strategy employed by Hu et al. (2023) accounted for both the error in the deterministic
798 predictions and the error in the learned perturbation distribution, effectively optimizing forecast accuracy and
799 ensemble spread at the same time. The authors assessed both the ensemble spread, and ensemble mean accuracy of
800 their model, and found that it had a better ensemble spread than simpler alternative ensemble generation strategies.
801 They also found that it had lower latitude-weighted RMSE than the ECMWF IFS to 5 days lead time for 2m
802 temperatures and total precipitation. ECMWF data beyond 5 days was not shown, but the SwinVRNN models had
803 latitude-weighted RMSE values lower than a weekly climatology baseline for three of the four variables shown to 14
804 days lead time. Bi et al. (2022) achieved a significant milestone with their model Pangu-Weather, the first ML-based
805 model to perform better than the ECMWF IFS to a lead time of 7 days based on RMSE and Anomaly Correlation
806 Coefficient (ACC) across several variables including geopotential height and temperature at 500 hPa. While they did
807 explore the utility of Pangu-Weather for ensemble generation, their approach was more simplistic than that
808 demonstrated by Hu et al. (2023). Pangu-Weather featured two major innovations over previous contributions to this
809 space:

- 810 1. It used 3D (latitude, longitude and height) input grids trained against 3D output grids. This enabled different
811 levels of the atmosphere to share information, which was not possible in FourCastNet in spite of predicting
812 variables on multiple atmospheric levels, because the levels were treated independently. In contrast, Pangu-
813 weather adopted a 3D convolutional method that the authors name the 3D Earth-specific transformer
814 (3DEST), which enabled the flow of information both horizontally and vertically.
- 815 2. It was made up of a series of models trained with different prediction time gaps. The motivation for this was
816 that, as noted by the authors, when the goal is to produce forecasts to 5 days (for example), but the timestep
817 of the basic forecast model is relatively short (e.g. 6 hours), many iterative executions of the model are

818 required, with the errors of each iteration feeding onto the next. A shorter model timestep results in greater
819 overall errors (due to more iterations being required to reach the final forecast lead time), and a longer model
820 timestep reduces this error. Motivated by this, the authors trained several versions of their model to predict
821 to different timesteps on a single iteration. The overall forecast to a given lead time was then constructed
822 using the longest possible timesteps. For example, for a 7-day forecast, a 24-hour forecast is iterated 7 times,
823 whereas for a 23-hour forecast, a 6-hour forecast is iterated 3 times, followed by a 3-hour forecast 1 time,
824 and 1-hour forecast 2 times. The authors noted that this strategy was not effective to multiweek or longer
825 timescales; they reported that training the model with a 28-day timestep was difficult, for example, and
826 suggested that more powerful or complex ML methods would be required to achieve this.

827 As well as the relatively broad measures of RMSE and ACC, the authors assessed the ability of their system to
828 represent the intensity and track of selected tropical cyclones. They found that Pangu-Weather predicted the tracks of
829 the cyclones considered with a high degree of accuracy compared to the ECMWF IFS, however it underestimated
830 cyclone intensity. The authors attributed this to the training data they used (ERA5) also underestimating cyclone
831 intensity. As noted above, the authors also explored the potential for producing useful ensemble forecasts. To assess
832 ensemble predictions, they perturbed the initial state of the system with Perlin noise vectors to produce a 100-member
833 ensemble of forecasts and calculated the RMSE and ACC of the ensemble mean for selected variables. As in Weyn et
834 al. (2021), the authors noted that the ensemble mean forecasts performed worse than a single deterministic forecast
835 for shorter lead times (e.g., 1 day), but better for longer lead times. Unfortunately, as with Pathak et al. (2022), Bi et
836 al. (2022) did not investigate the properties of the spread of the ensemble or assess its skill using standard probabilistic
837 skill metrics, and their approach to ensemble generation was much simpler than that of Hu et al. (2023).

838 As already mentioned above, the skill of Pangu-Weather was exceeded by GraphCast, although Lam et al. (2022) only
839 assessed GraphCast in a deterministic setting. Nonetheless, there is nothing stopping GraphCast from being used to
840 generate ensemble forecasts in a manner similar to Pangu-Weather. The authors of this review look forward to a more
841 in-depth intercomparison of the pure ML models in the literature, including an assessment of their performance for
842 ensemble predictions.

843 Although the ensemble systems presented in Weyn et al. (2021) and Hu et al. (2023) had lower overall accuracy than
844 the other models discussed in this section, they still represented the most comprehensive analysis of the behavior and
845 performance of ensemble ML models (in terms of considering optimal ensemble perturbation strategies, and
846 quantifying the ensemble behavior) at the time of writing this review. Further investigation into the best methods to
847 generate and evaluate pure ML model ensembles would be a highly beneficial contribution to the field.

848 **5.4. Moving to more extensible models**

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

853 This need for greater flexibility in both the input data sources and predictive outputs of ML weather and climate
854 models was recognized by Nguyen et al. (2023), who developed a transformer architecture-based ML model called
855 ClimaX. This model was designed as a foundational model, trained initially on datasets derived from the CMIP6
856 (Eyring et al., 2016) dataset, and able to be readily retrained to specific tasks using transfer learning. The authors
857 demonstrated the skill of ClimaX against simpler ML models, and in some cases a numerical model (ECMWF IFS),
858 for a variety of tasks including weather prediction, sub-seasonal prediction, climate scenario prediction, and climate
859 downscaling. The authors showed that ClimaX was able to make skillful predictions in scenarios unseen during the
860 initial CMIP6 training phase. Furthermore, ClimaX used novel encoding and aggregation blocks in its architecture to
861 enable greater flexibility in the types of variables used for training, and to reduce training costs when a large number of
862 different input variables were used.

863 **5.5. Benchmark datasets for ML weather models**

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

866 As the interest and activity in the use of ML as a potential alternative to knowledge-based numerical GCMs has grown,
867 the need for consistent benchmarks for the intercomparison of ML-based models has become increasingly clear. Rasp
868 et al. (2020) addressed this need with the introduction of WeatherBench. On this platform, the authors provided data
869 derived from the ERA5 archive that has been simplified and streamlined for common ML use cases and use by a broad
870 audience. They also proposed a set of evaluation metrics which facilitate direct comparison between different ML
871 approaches, and provided baseline scores in these metrics for simple techniques such as linear regression, some deep
872 learning models and some GCMs. Since the publication of WeatherBench, more benchmark datasets tailored to other
873 domains have been created, including RainBench (de Witt et al., 2020), WeatherBench Probability (Garg et al., 2022),
874 and ClimateBench (Watson-Parris et al., 2022). Weyn et al. (2020) chose datasets and assessment metrics consistent
875 with WeatherBench to facilitate intercomparison of results. Rasp & Thuerey (2021) directly used the benchmarks
876 provided by WeatherBench in their assessment. They demonstrated that their model outperformed previous
877 submissions to WeatherBench, highlighting its value as a tool to allow intercomparability of ML-based weather
878 models. Other examples of studies using WeatherBench data and analysis methods are Clare et al. (2021) and Weyn
879 et al. (2021). The parameters of a good benchmark dataset were further elucidated by Dueben et al. (2022), who
880 provided an overview of the current status of benchmark datasets for ML in weather and climate in use in the research
881 community and provided a set of guidelines for how researchers could build their own benchmark datasets.

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

888 It should also be noted that all of the major milestones and high-profile ML models described in this section so far
889 have relied to some degree or another on reanalysis datasets produced by physics-based models. The provision of
890 higher resolution and higher quality open datasets have the potential to drive progress in this area as much as, if not
891 more than, improvements and further research into ML algorithms.

892 **5.6. A hybrid approach**

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

900 **5.7. ML for predicting ocean variables**

901 More recently, greater attention has been paid to the application of ML to the ocean, particularly for seasonal to multi-
902 year prediction. Initial work in this space focused on directly predicting key indices such as the NINO 3.4 index. For
903 example, Ham et al. (2019) trained a CNN to produce skillful El Niño Southern Oscillation (ENSO) forecasts with a
904 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
905 availability of observational data for training. To overcome this, the authors used transfer learning[†] to train their model
906 first on historical simulations, and then on a reanalysis from 1871 to 1973. Data from 1984 to 2017 was reserved for
907 validation. Ham et al. (2021) improved on this by including information about the current season in the network inputs
908 as one-hot vectors[†]. Including this seasonality information led to an overall increase in skill relative to the model in
909 Ham et al. (2019), in particular for forecasts initiated in boreal spring, a season which is particularly difficult to predict
910 beyond.

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

918 Taylor & Feng (2022) moved from prediction of indices to spatial outputs, training a Unet-LSTM[†] model on ECMWF
919 ERA5 monthly mean Sea Surface Temperature (SST) and 2-m air temperature data from 1950-2021 to predict global
920 2D SSTs up to a 24-month lead time. The authors found that their model was skillful in predicting the 2019-2020 El
921 Niño and the 2016-2017 and 2017-2018 La Niñas, but not for the 2015-2016 extreme El Niño. Since they did not

922 include any subsurface information in their training data (in contrast to Ham et al. (2019) and Ham et al. (2021), who
923 included ocean heat content), they concluded that subsurface information may have been relevant for the evolution of
924 that event.

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

928 **5.8. ML for climate prediction**

929 The literature on the use of ML for prediction on seasonal to climate timescales is still relatively sparse compared to
930 its use for nowcasting and weather prediction. Some examples have been covered in previous sections, such as Weyn
931 et al. (2021) on subseasonal to seasonal timescales in the atmosphere, and Ham et al. (2019), Ham et al. (2021), Kim
932 et al. (2022) and Taylor & Feng (2022) on seasonal to multiyear timescales in the ocean. A major cause for this sparsity
933 is that deep learning typically requires large training datasets, and the available observation period for the earth system
934 is too short to provide appropriate training data for seasonal to climate timescales in most applications. On the
935 subseasonal to seasonal end, this may be overcome by including more slowly-varying fields in the training (e.g. ocean
936 variables), by designing models to learn the underlying dynamics which drive long-term variability, and by including
937 more physical constraints on the models. On the climate end these same methods could be beneficial, as well as
938 transfer learning, as is done in Ham et al. (2019), and data augmentation[†] techniques. Additionally, interest is
939 increasing in the use of ML to predict weather regimes and large-scale circulation patterns, which may prove beneficial
940 in informing seasonal and climate predictions (Nielsen et al., 2022). Watson-Parris (2021) argued that the differences
941 between NWP to multiyear prediction and climate modelling mean that the ML approaches best suited to each can be
942 very different. This may also help to explain why the rapid pace of advances in ML based weather models has not
943 translated into a similar trend in climate modelling.

944 Despite this, with the growing maturity of the field of ML for weather and climate prediction, there is every reason to
945 believe the challenges of prediction on seasonal to climate timescales can be overcome.

946 **6. Physics constrained ML models**

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

952 An excellent review of the possible methods for incorporating physics constraints into ML models for weather and
953 climate modelling, along with 10 case studies of noteworthy applications of these methods, is presented in Kashinath
954 et al. (2021). The scope of Kashinath et al. (2021) is broad and includes studies not applied directly in the context of

955 weather and climate modelling, but applicable to it. Rather than repeat the total of this summary here, the reader is
956 directed to this review.

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

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

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

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

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

976 **7.1. Nudging**

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

985 **7.2. Uncertainty quantification**

986 A common criticism of some ML models such as NNs is that it is difficult to represent the uncertainty of their outputs.
987 Some examples of studies that have sought to overcome this have already been mentioned in Section 3.8, and there

988 are other examples in the literature (e.g. Grigo & Koutsourelakis, 2019; Atkinson, 2020; Yeo et al., 2021; O'Leary et
989 al., 2022), however it is nonetheless still a relatively underexplored aspect of ML models for physical systems. Psaros
990 et al. (2022) suggest that this may be because they are also under-utilized within the broader deep learning community,
991 and it is thus a developing field that is not universally trusted and understood yet. They also point out that the physical
992 considerations inherent to ML applied to physical systems often make them more complicated and computationally
993 expensive than standard ML applications, further disincentivizing the inclusion of uncertainty quantification in an
994 already complex problem.

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

999 A related topic which is facing similar challenges is the question of explainability of ML approaches; often there is
1000 value in understanding the relative roles and importance of predictors in an ML model, or the relative significance of
1001 different regions of the predictor data. Flora et al. (2022) provide a good overview of approaches to this and compare
1002 their relative drawbacks and benefits.

1003 **7.3. Capturing extremes**

1004 While there is now an abundance of examples of ML being used for model parameterization schemes, full model
1005 replacement, downscaling, and PDE solvers (much of which is covered in this review), there are relatively few
1006 examples which address the question of how well ML approaches can reproduce extreme events and statistics, both
1007 in terms of the distribution of values predicted in a single-member (i.e., non-ensemble and non-probabilistic) ML
1008 model and in terms of the distribution of predicted outcomes in a probabilistic or ensemble ML model.

1009 Both Pathak et al. (2022) and Bi et al. (2022), introduced in Section 5.2, investigated the ability of their models to
1010 correctly represent extremes, using a similar approach. They divided their test dataset into 50 percentile bins
1011 (distributed logarithmically by Pathak et al. (2022) and linearly by Bi et al. (2022)) between the 90th and 99.99th
1012 percentiles, and computed the relative quantile error between their forecast and ground-truth as a function of lead-
1013 time. Pathak et al. (2022) note that they set their highest percentile bin at 99.99% because of the small sample of
1014 datapoints beyond this percentile making a statistically significant analysis difficult. Both Pathak et al. (2022) and Bi
1015 et al. (2022) found that their models consistently under-forecast extremes to a greater degree than the ECMWF IFS.

1016 Watson (2022) presents a strong argument for the need for a greater focus on the ability of ML weather and climate
1017 models to be able to predict extremes in order for them to meet the needs of users. They present a summary of some
1018 examples of ML models which have sought to predict extreme events according to certain return period definitions.
1019 The example most relevant for this review is Lopez-Gomez et al. (2023), who used a NN with a custom loss function
1020 that preferentially weighted extremes to predict global extreme heat. They found that their custom loss function led to
1021 improved representation of the tails of the distribution (i.e., predictions of extreme heat), and, interestingly, did not
1022 result in any major loss of performance for the middle of the distribution.

1023 The under-prediction of extremes seen in Pathak et al. (2022) and Bi et al. (2022) is consistent with the findings of
1024 Lopez-Gomez et al. (2023), given that neither were not optimized for predicting extremes. These findings all point to
1025 the idea that in order for ML weather and climate models to be able to skillfully predict extreme events, model training
1026 regimes, loss functions and architectures will need to be employed which take into consideration ways to optimize for
1027 these regimes.

1028 **7.4. Object identification within models**

1029 An alternative to achieving greater model accuracy and skill for predicting extremes through increasing resolution of
1030 the entire model grid is to develop techniques to identify critical systems and physical phenomena within the model,
1031 and embed higher resolution temporary subgrids or specialized models within the larger GCM to more accurately
1032 simulate those processes. A challenge to overcome to achieve this is automatically identifying key model features,
1033 since it typically requires a labelled dataset. This requirement can however be avoided, and a variety of both supervised
1034 and unsupervised machine learning approaches to object detection have been demonstrated in the literature.

1035 Mudigonda et al. (2017) were a relatively early example of the application of ML to this challenge. They investigated
1036 the feasibility of using a variety of NN architectures to identify storms, tropical cyclones and atmospheric rivers within
1037 model data, with promising results. Prabhat et al. (2021) provided a valuable resource to the community with their
1038 development of ClimateNet, a labelled open dataset and ML model for the segmentation and identification of tropical
1039 cyclones and atmospheric rivers. This was used by Kapp-Schwoerer et al. (2020) to train a NN to identify and track
1040 these extreme events in Community Atmosphere Model 5 (CAM5; Conley et al. 2012) data. O'Brien et al. (2021)
1041 considered the need for uncertainty quantification in object identification, using a Bayesian approach to build an
1042 atmospheric river detection framework. Finally, Rupe et al. (2023) took a physics-informed approach to object
1043 detection, defining 'local causal states' using speed-of-light causality arguments to identify regions of organized
1044 coherent flow and bypassing the requirement for labelled datasets. They demonstrated the utility of their approach for
1045 the unsupervised identification and tracking of hurricanes and other examples of extreme weather events.

1046 While there are unsupervised learning approaches which have shown value for object detection in weather and climate
1047 data (e.g. Rupe et al., 2023), a major limitation of this area of research is the shortage of labelled datasets for supervised
1048 learning methods , with ClimateNet being an isolated example.

1049 **7.5. GPUs and specialized compute resources**

1050 GPUs and TPUs are specialized hardware which are well suited to highly parallelizable matrix operations, ideal for
1051 solving neural network operations. TPUs have been developed specifically for deep learning applications. Both GPUs
1052 and TPUs are likely to be available on many of the next generation of supercomputers, but much of the current Fortran-
1053 based numerical weather and climate model infrastructure cannot be run on them in their current state. Data
1054 bottlenecks also exist between the GPUs (which have their own on-board memory) and the main memory accessible
1055 to the CPU. While efforts are underway to make numerical and climate models better suited to GPUs, for example
1056 with the development of LFRic (Adams et al. 2019), the new weather and climate modelling system being developed

1057 by the UK Met Office to replace the existing Unified Model (Walters et al. 2017), there is still a long way to go before
1058 entire weather and climate models can be reliably run on GPU or other specialized compute architectures. At the same
1059 time, some neural network designs are aimed squarely at the partial differential equation solving at the core of
1060 numerical methods. Since neural network evaluation utilizes simpler mathematical operations than current PDE
1061 solvers, they offer the prospect of significant computational advantages on non-specialized (i.e., CPU) hardware.

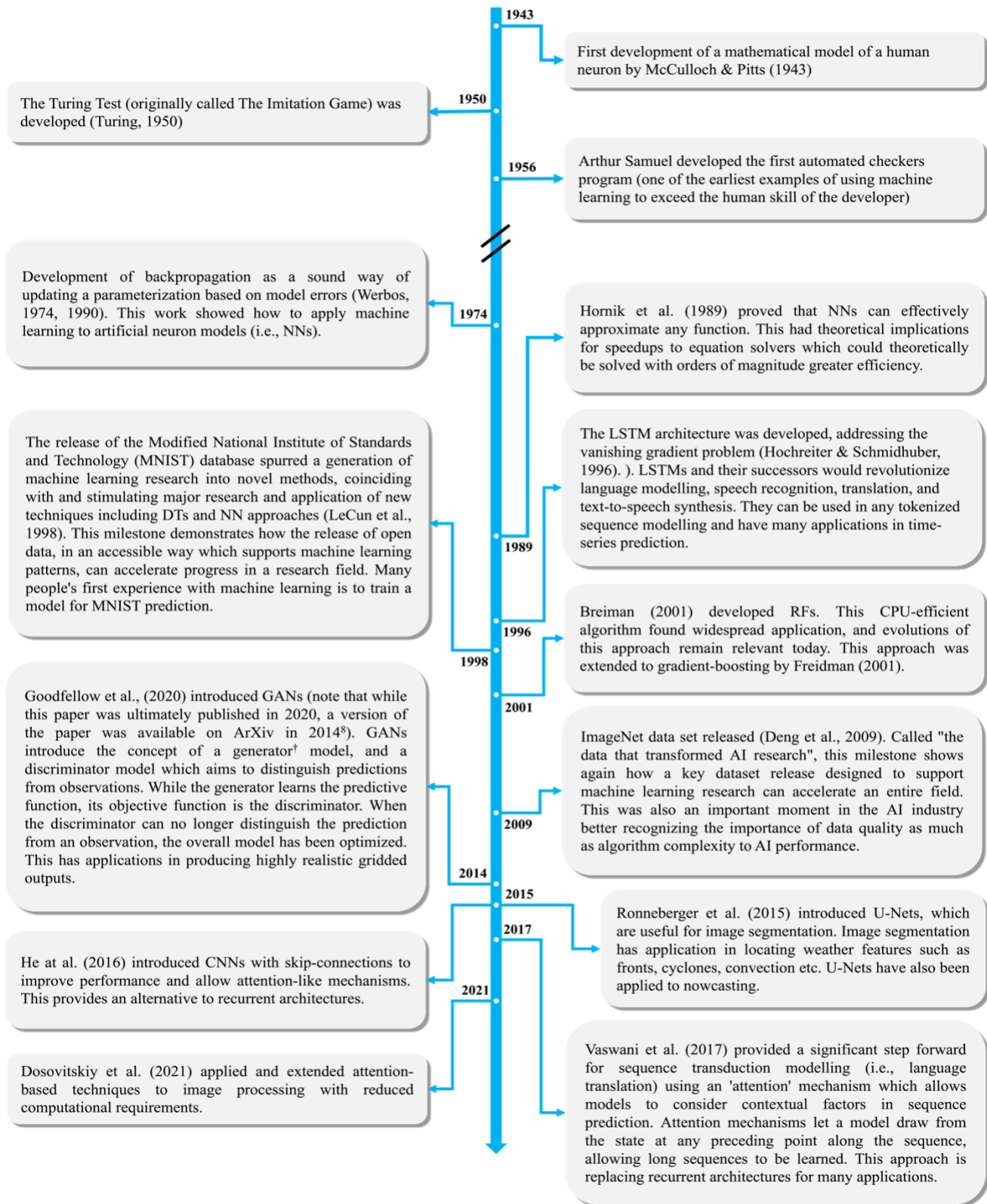
1062 **8. Perspectives on machine learning from computer science**

1063 This section provides a brief perspective on weather and climate modelling from the computer science domain and
1064 aims to provide the earth system scientist with a short list of the main relevant innovations in computer science. As
1065 was noted in Section 1, ML models are often regarded as black-boxes, largely because of the design of many prominent
1066 ML systems. In principle, it is not quite right to refer to the trained model as "a machine learning model", in the sense
1067 that the process of training the model is "machine learning", once the model is trained it is definable by a set of
1068 mathematical equations and coefficients, much like any physical, statistical, or theoretical model. Thus the machine
1069 learning refers to the training process, not the model itself. The essence of ML is the level of automation involved.
1070 Even in typical ML models such as large NNs, the model architecture is typically specified manually by the data
1071 scientist or physical scientist involved. The automated derivation of model architecture and composition is not yet
1072 mature for large models, although it is explored through evolutionary programming techniques whereby the learning
1073 of architecture as well as parameterization is automated.

1074 The complex nature of the Earth system means that ML models which seek to emulate it (or subcomponents of it) will
1075 likely also need to be quite complex, and will contain a mixture of ML architectures and algorithms. This is borne out
1076 by the increasing degree of complexity and variety seen in the ML models in the literature reviewed in previous
1077 sections.

1078 A large degree of the current research focus is on very large or deep NNs which rely both on the universal
1079 approximation theorem and practical experimentation to capture a prediction function without needing to explicitly
1080 represent the processes being modeled. In a conceptually similar fashion to how a Fourier decomposition can represent
1081 any wavelike function, the universal approximation theorem establishes that a NN may approximate any function,
1082 subject to its size and the required degree of accuracy (Hornik, Stinchcome and White 1989). Deep learning has been
1083 highly effective in approaching many problems, but many limitations are acknowledged, as evidenced by the current
1084 widespread focus on trustworthy computing and efforts towards explainable ML systems. Some ML models take a
1085 direct approach to modelling the uncertainty of the system being simulated by representing the model state variables
1086 as a probability distribution or degree of confidence. Many contemporary weather and climate model derive their
1087 probabilistic outputs from an ensemble of perturbed members, however an alternative approach is to represent each
1088 part of the belief state[†] of the model as a distribution or likelihood, built up either empirically or by fitting a gaussian
1089 or other known distribution (e.g., Clare et al., 2021).

1090 A timeline of some key innovations in ML is presented in Figure 4. The scale of the timeline is broken between 1956



1091

1092 Figure 4: A timeline of key breakthroughs in ML.

1093

1094 and 1974, and Taking that gap in progress into account, it is clear from this visualization that the rate of innovation in
1095 ML has increased significantly over the last 35 or so years. This is likely driven by a range of factors including the
1096 increasing availability of compute resources suited to ML applications, and the explosion of available data for training.
1097 This history shows the degree and rate of research into processing images, text and other sequences based on semantic
1098 understanding of content, but does not demonstrate capturing physical processes as a core element. Advances in the
1099 weather and climate modelling domain have a more explicit goal of properly portraying real physical processes.
1100 Bringing these concepts together promises to uplift capability in both fields.

1101 **9. Practical Perspectives on Machine Learning for Weather and Climate Models**

1102 A major driver of research into, and improvement of, weather and climate models is increasing the skill of operational
1103 forecast systems, and increasing the accuracy and trustworthiness of climate projections. Therefore, an important
1104 consideration for ML in the context of weather and climate models is the need for it to ultimately be integrated into a
1105 complete predictive system with practical application for forecasting or climate projections.

1106 However, the research findings covered in this review, in spite of being compelling, are yet to make major changes to
1107 operational modelling systems, or standard climate projections.

1108 We have identified three major challenges facing the transition of ML-based innovations into operational settings.
1109 Similar challenges are faced in the context of climate projections, however since these are out of scope for this review
1110 we do not discuss them directly, and instead leave them as a topic for other publications.

1111 The first challenge is the need to assess when a research finding is sufficiently compelling and robust to justify
1112 integration into established operational systems. Since the major function of operational meteorological services is to
1113 inform of future conditions, largely for managing risk or optimizing benefits, a conservative approach is taken to
1114 changing these systems. The utmost premium is put on accuracy, resilience, reliability, and solid scientific foundation,
1115 and many novel research finding require extensive further evaluation and development before they can be considered
1116 ready for inclusion into operational systems. Understanding when to invest this degree of effort in bringing a research
1117 innovation into a major model or scientific configuration upgrade can be difficult.

1118 The second major challenge is establishing the right balance between potentially unwieldy monolithic ML models
1119 which predict all variables of interest, and many smaller limited scope models which each focus on predicting one or
1120 a small number of variables well. The former option is more similar to current dynamical systems, while the latter
1121 option is potentially more easily achievable using an ML approach, but risks becoming difficult to manage due to the
1122 proliferation of small, separate systems. The early effectiveness of limited-purpose ML models provides the ability to
1123 augment existing services without disruption, however aside from the logistical complexity of many small systems, a
1124 risk associated with this approach is that inconsistencies between predictions may arise from their independent
1125 forecasts, leading to confusion from users and an erosion of trust.

1126 Finally, the third major challenge is how to best monitor and maintain the skill of ML-based systems in a real-time
1127 operational context. Explainability of ML systems is an emerging field, and is not yet sufficiently mature for
1128 application to real-time operational monitoring. Until this changes, the ongoing trustworthiness of operational ML

1129 systems will be difficult to demonstrate. Similarly, online learning in ML weather and climate models is not yet a well
1130 explored research area. The use of online learning is likely to be important for operational ML models to be able to
1131 develop resiliency and maintain good skill over time, so more work will be needed in this area before these models
1132 can see greater uptake in operational systems.

1133 In addition to these major challenges, agencies looking to incorporate ML components into their operational systems
1134 must consider that:

- 1135 • the explainability of ML model errors in the case of poor forecasts that may come under scrutiny,
- 1136 • the robustness of ML models to real-time data issues such as data dropouts or input data degradation must be
1137 established, and
- 1138 • the lack of infrastructure in these agencies to support ML models in an operational setting will need to be
1139 addressed.

1140 Operational development is typically quite incremental, and it is likely that progress will be made in small achievable
1141 steps along the evolving technical frontier. However promising and fascinating as a research direction, full model
1142 replacement with ML alternatives is currently not mature enough for an operational setting. Instead, the authors predict
1143 that the first types of ML systems to be seen in operations will include parameterization scheme replacements and
1144 emulators, solver replacements, super-resolution, new approaches to data assimilation of novel observation sources,
1145 and both pre- and post-processing applications (although of course not all of these have been covered in this review).
1146 It is expected that the research into, and application of, ML methods will represent a growing proportion of weather
1147 and climate model research, with increasingly sophisticated and skillful model components finding their way into
1148 major model releases over the coming years. These components are appealing for both computational and model skill
1149 reasons, and are expected to be highly promising avenues of research.

1150 **10. Ethical considerations for Machine Learning for Weather and Climate Models**

1151 Not all papers in this review included a discussion of the ethical considerations associated with using machine learning,
1152 nor necessarily touched on what constitutes a sufficiently rigorous verification methodology for machine learning
1153 models. There is a clear relationship between ethical considerations, the explainability of models, and the rigor of
1154 verification applied to ensure that models behave as expected under a variety of conditions (and do not include
1155 unexpected behaviours).

1156 While this review paper does not provide an introduction to AI and ML ethics in general, a brief overview of some
1157 of the important considerations for the application of ML in the context of weather and climate modelling is
1158 provided in this section. Ethical frameworks vary in different cultural and geographical contexts, and for a more
1159 general introduction to the ethical considerations surrounding AI and ML, the reader is directed to the paper
1160 *Recommendations on the Ethics of Artificial Intelligence* (United Nations Educational, Scientific and Cultural
1161 Organisation (UNESCO), 2022).

1162 For ML applied to weather and climate modelling, some considerations to ensure sufficient robustness and reliability
1163 include whether:

- 1164 • testing, training and validation data sets are sufficiently representative of the data in general
 - 1165 • potential causal correlations between testing, training and validation data have been treated correctly
 - 1166 • trained models have been tested for reliability against adversarial examples
 - 1167 • data augmentation (e.g. noise addition) has been utilized to enhance model robustness
 - 1168 • an evaluation of the potential for model drift has been performed
 - 1169 • the training data is biased in a way which results in ethical unfairness (for example – remote communities
 - 1170 may not receive equal-skill predictions due to a lack of observational training data in remote areas,
 - 1171 • the machine learning method is compared to a suitable alternative, such as a known physical model in
 - 1172 addition to any comparisons to machine learning models or the provision of aggregate statistics
 - 1173 • the data that has been used has been gathered ethically, and any personal information has been treated
 - 1174 properly (such as when processing weather reports from individuals)
 - 1175 • the authors have identified any caveats regarding ethics, reliability, robustness or explainability
 - 1176 • the authors have investigated the physical realism of the predictions from ML models
- 1177 This list is not comprehensive, however. A thorough overview of the explainability, reliability, ethics, and
- 1178 verification of ML models in weather and climate has not been covered in prior literature and the field will
- 1179 benefit from further work in this area.

1180 **11. Future research directions**

1181 The already-demonstrated and potential future applications for ML in weather and climate modelling are significant

1182 in number, and identifying the most fruitful avenues for future research can seem overwhelming. A good

1183 understanding of the current state of the weather and climate modelling field, along with knowledge of the key

1184 developments in ML research, are required to assess the potential benefits of a given research direction.

1185 As can be seen from the timeline of machine learning presented in Figure 4, older techniques can prove to be

1186 relevant many years later, and there are many techniques from computer science which may become relevant for

1187 contemporary weather and climate modelling problems and research.

1188 Furthermore, due to the general applicability of many ML approaches, research progresses in one subdomain may

1189 have implications and benefits for another. For example, DeepONets were developed for, and shown to be

1190 successful for, solving PDEs, but were adopted by Pathak et al. (2022) for their pure ML model FourCastNet with

1191 great success.

1192 To help the reader navigate the myriad research areas where ML for weather and climate modelling could be

1193 progressed, five categories of future research directions are presented in Figure 5, along with some specific areas of

1194 research, and benefits that could arise from them.

1195 These categories are not mutually exclusive – indeed there is overlap between the research areas and benefits

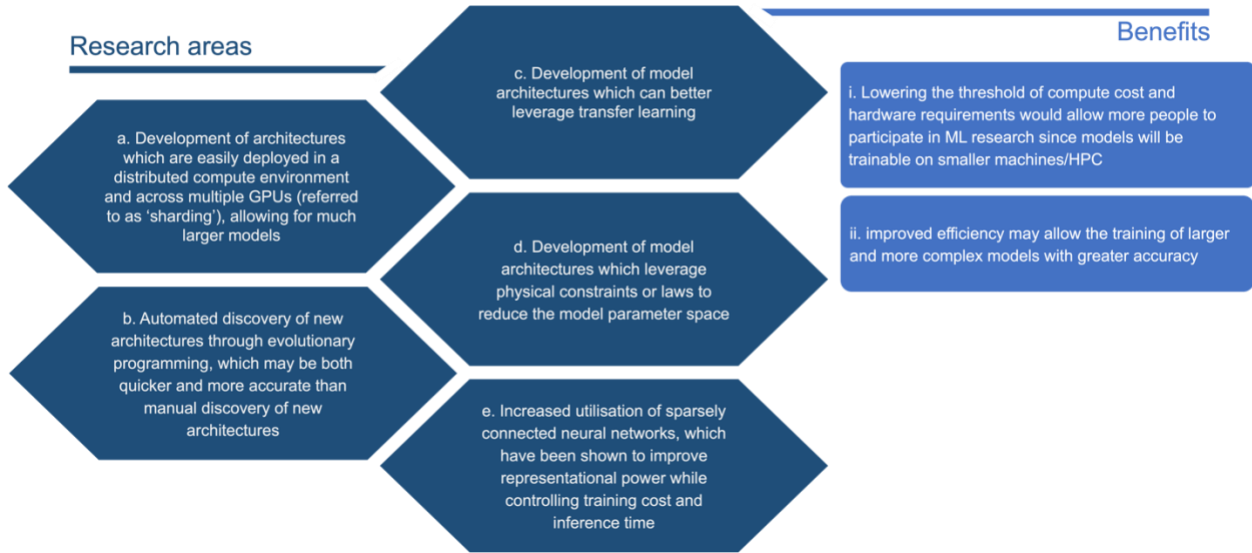
1196 highlighted in each category (for example, some research foci in Categories 2 and 3 are also applicable to Category

1197 5). The groupings are instead intended to help guide the focus of researchers, and to provide a quick overview of the

1198 key topics where the community would most benefit from research progress.

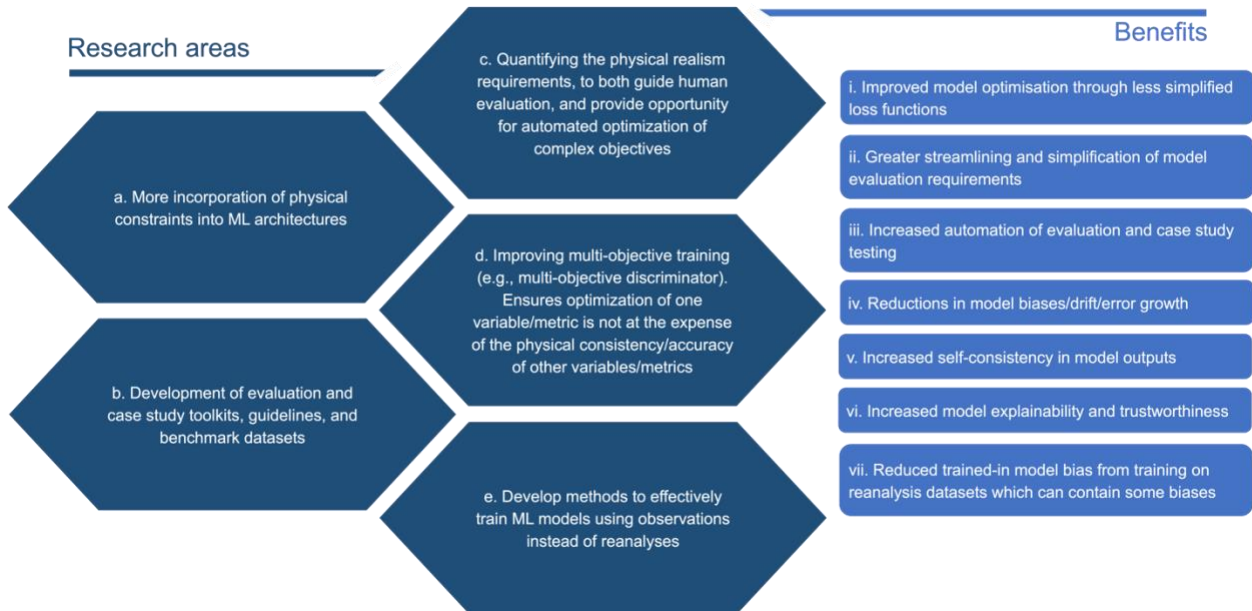
1199
1200

Category 1: Improving training speed and efficiency



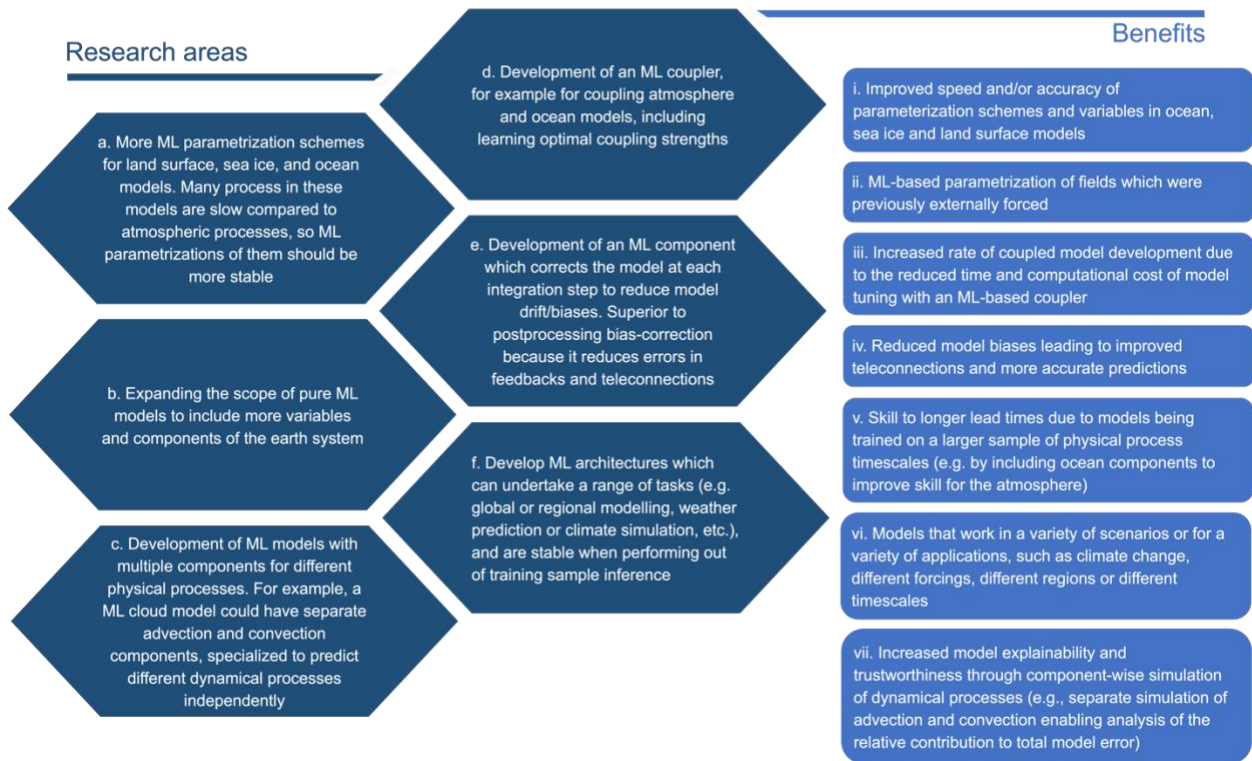
1201

Category 2: Physically consistent/constrained models and evaluation



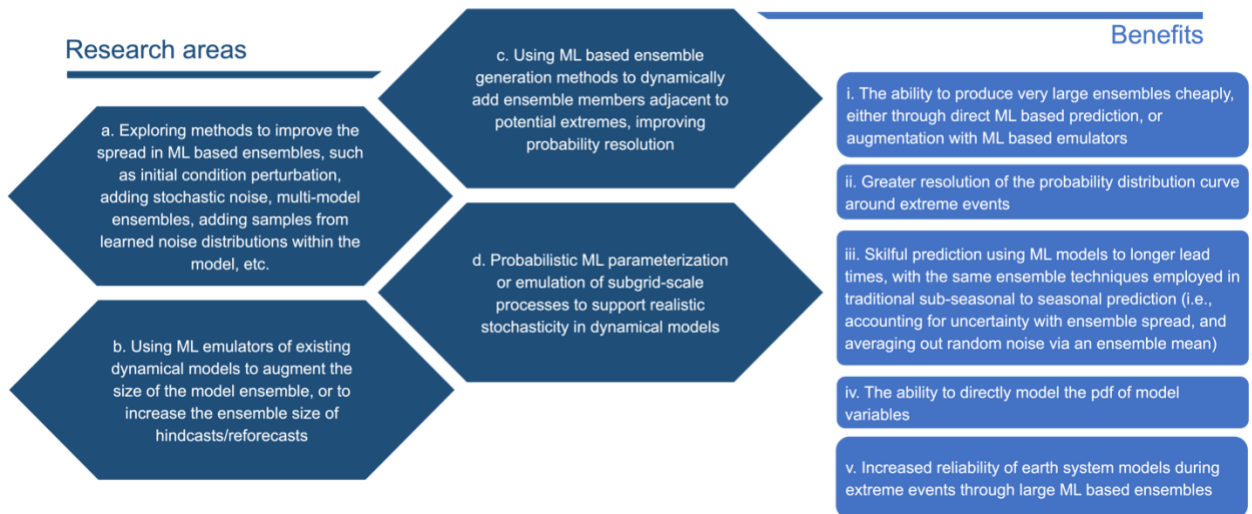
1202

Category 3: Weather and climate modelling domain specific research



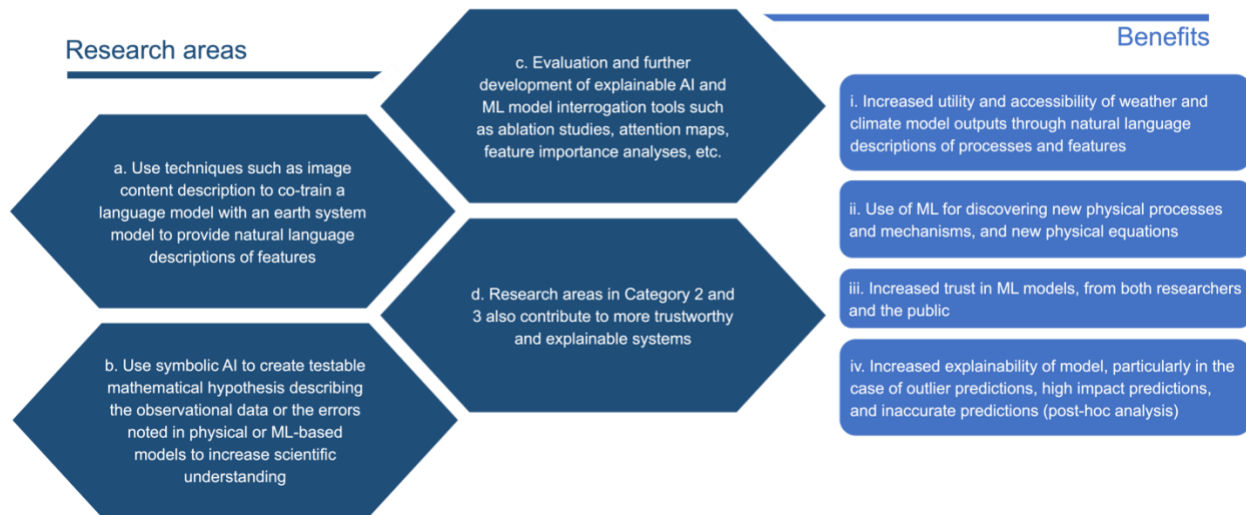
1203

Category 4: Probabilistic prediction



1204

Category 5: Trustworthy and explainable systems



1205
1206 Figure 5: Five categories for future ML research, including suggested research focusses for the community in each
1207 category, and potential benefits which could be realized by research and development progress.

1208
1209 Many of the research areas presented are complementary to each other, for example progress in making ML models
1210 more affordable to train (Category 1) will increase the utility of ML solutions to a wider community of researchers,
1211 and will likely accelerate the rate of progress in the other categories. Progress in the use of physically-informed
1212 approaches (e.g. Category 2, area a., or Category 3, area c.) could also lower the training cost of models by reducing
1213 the degree of redundancy in the model. On the other hand, approaches such as Category 3, area f., leading to an
1214 outcome such as benefit vi. would potentially reduce the demand for more cheaply trainable models, since they
1215 could be readily turned to a variety of tasks, saving researchers the need to train their own model from scratch.
1216 The research areas and ideas presented here are by no means a comprehensive list. Rather they are intended to be
1217 used as a source of inspiration, and the authors of this review are excited to see where the community chooses to
1218 focus their efforts in the coming years.

1220 12. Conclusions

1221 In this review we have presented a comprehensive survey of the literature on the use of ML in weather and climate
1222 modelling.

1223 We have found that the ML models being most often explored include RFs and NNs, with a high prevalence of FCNNs
1224 and CNNs. We have also identified some recent innovations which have proven to be highly effective in the weather
1225 and climate modelling space, including DeepONets and variants thereof, Graph NNs, and PINNs.

1226 This review has demonstrated that ML is being successfully applied to many aspects of weather and climate modelling.
1227 We have presented examples from the literature of its application in (1) the emulation and replacement of subgrid-
1228 scale parametrizations and super-parametrizations, (2) preconditioning and solving of resolved equations, (3) full
1229 model replacement, and (4) a selection of other adjacent areas.

1230 Nonetheless, there are still many research challenges to overcome, including:

- 1231 • addressing the instabilities excited in physical models due to the inclusion of ML components;
- 1232 • increasing the ease of technical integration (in particular, Fortran compatibility);
- 1233 • memory and computational concerns;
- 1234 • representing a sufficient number of physical parameters and increasing physical and temporal resolution in
- 1235 ML-based weather and climate model implementations (which currently feature reduced fields and levels
- 1236 compared to physics-based numerical models);
- 1237 • moving from a focus on individual parts of the earth system (i.e., the atmosphere, the ocean, the land surface
- 1238 etc.) to tackling the challenges associated with coupled models (i.e., where models of individual components
- 1239 of the earth system are coupled together). Increasingly, operational weather and climate models are coupled
- 1240 land-atmosphere-ocean-sea-ice models in order to more accurately represent the relevant timescales and
- 1241 processes in the earth system, and ML modelling efforts need to reflect this;
- 1242 • more thorough evaluation of the physical realism of ML-based predictions, at various length-scales, across
- 1243 parameters, and looking at the three-dimensional structures
- 1244 • Exploring the use of generalized discriminators to augment traditional loss functions in model training (to
- 1245 achieve a multivariate generalized objective function)
- 1246 • the need for more good quality training data; and
- 1247 • the practical challenges of integrating ML components or models into an operational setting.

1248 This list, together with Section 11, provides a set of focus areas for future research efforts.

1249 If the current trend in skill gains in full ML weather and climate models continues, it is possible they will eventually
1250 be considered viable alternatives to traditional numerical models. However, in the meantime it is likely that ML
1251 components will replace an increasing number of physics-based model components, with models the near-term future
1252 being hybrid ML-physical models. A likely future scenario is one where the best weather and climate models are a
1253 blend of ML and physics-based components, deriving skill from both data driven and physical methodologies.

1254 Some possible avenues through which increases in ML-based weather and climate model skill might be achieved is
1255 by operating at higher resolutions, resolving more processes which are implicit in the training data, or by undertaking
1256 experiments on synthetic data to address the paucity of real-world data.

1257 Another benefit of ML approaches to weather and climate modeling is the relative computational cheapness of ML
1258 alternatives to current physics-based modelling systems. This has the potential to open the door to experiments that
1259 would not be feasible otherwise. For example, experiments requiring a very large ensemble would be more feasible
1260 with a computationally cheap ML approach.

1261 The literature reviewed here indicates that 'out of the box' ML approaches and architectures are not effective when
1262 used in a weather and climate modelling context. Rather, ML architectures must be adapted to satisfy conservation of
1263 energy, represent physically realistic predictions and processes, and maintain good model stability. At the same time,
1264 computational and memory tractability must be maintained.

1265 Advances in the sophistication, complexity and efficiency of ML architectures are being heavily invested in for many
1266 use cases in other disciplines and in the private sector (e.g., condition-action pose estimation, text to video generation,
1267 stable diffusion/text to image, chatbots, facial recognition, semantic image decomposition, etc.). In order to capture
1268 the full benefits of ML for the weather and climate modelling domain, academic and operational agencies will need
1269 to continue to support research in this space. This includes contributing to the research effort through foci such as
1270 those highlighted in Section 11 and in this section, and through addressing the particular challenges facing agencies
1271 interested in the operational and/or realtime deployment of ML based models as the basis for services or the provision
1272 of advice (discussed in Section 9).

1273
1274 Interest and progress in the application of ML to weather and climate modelling has been present for close to 30 years,
1275 and has begun to accelerate rapidly in the last few years. There is good reason to believe that ML as a tool will have
1276 transformational benefits and offers great potential for further application in weather and climate modelling.

1277 **Machine Learning Glossary of Terms**

1278 This glossary includes terms which the reader will come across frequently in machine learning literature for the
1279 weather and climate, as well as in machine learning literature generally. Most of these terms are used in this paper
1280 while others support further reading.

1281 **Activation Function.** The function which produces a neuron's outputs given its inputs. Commonly, this includes a
1282 learned bias term which is added to the data inputs before evaluation with a single function to produce the output
1283 value. Examples of the functions used include linear, sigmoid and tanh.

1284 **Adversarial attack.** The deliberate use of malicious data input in a real-world setting intended to cause a
1285 misclassification, underperformance or unexpected behaviours. Examples include emails designed to avoid spam
1286 filters, or images that have been modified to avoid recognition.

1287 **Adversarial example.** A specialised input which results in a misclassification or underperformance of a predictive
1288 model. An example of this concept is an image which has had subtle noise added to it resulting in a copy of that image
1289 which is visually indistinguishable from the original, but which nonetheless causes a misclassification. The term
1290 'adversarial' is used to refer to the way the example fools the model and is not necessarily intended to convey the
1291 sense of malicious intent, although the term is often applied in that fashion. Adversarial examples demonstrate that
1292 machine learning models may be more brittle than expected based on ordinary training data alone. To increase model
1293 robustness, adversarial examples may be generated and added to the training set. Data augmentation techniques such
1294 as flipping, warping and adding noise (any many other techniques) are also used to generate additional training data
1295 to increase robustness and performance.

1296 **Attention mechanism.** A mechanism to allow sequence prediction models to increase the importance of key terms
1297 within that sequence which may be nonlocal and modified in meaning according to the other terms of the sequence.

1298 **API.** Application Programming Interface. A set of programming functions, methods or protocols by which to build
1299 and integrate applications. APIs may be "web" APIs or imported from software packages in which case they are more
1300 often referred to as libraries.

1301 **Autoencoder.** A neural network architecture which learns to produce a 'code' for an input sequence from which the
1302 original data can be retrieved. The code is shorter than the original input sequence. Applications include data
1303 compression and denoising data.

1304 **Back propagation.** A process of utilising the errors from a prediction to update the weights and biases of a neural
1305 network.

1306 **Batch.** See training batch.

1307 **Batch normalisation.** Data normalisation which aligns the means and variances of input data to a model. For
1308 computational reasons, this is performed separately for each training batch.

1309 **Belief state.** The current state of the world which is believed to be true according to a model. A common architecture
1310 in realtime applications whereby a belief state is updated according to an update function on the basis of new
1311 observations.

1312 **Channel.** An additional dimension to data which is usually not a spatial dimension. Examples include the red, green
1313 and blue intensity images which comprise a colour image. Another example could be to represent both temperature
1314 and wind speed as channels.

1315 **Classification.** A model which attempts to diagnose or predict the category, label, class or type that an example falls
1316 within.

1317 **Climatology.** Refers to the usual past conditions for a location at a time of year. Usually calculated by temporal mean
1318 across years of a dataset, for a given time interval within those years (e.g., for a dataset of monthly mean values
1319 spanning all months of all years from 1990 to 2020, the monthly mean climatology would be obtained by averaging
1320 across all the Januarys from each year, all the Februarys, etc., to obtain an "average January", an "average February",
1321 etc.). Climatologies are often used in the same manner as persistence as a baseline prediction against which to measure
1322 a predictive model. For example, a model predicting a value for January could be compared to the climatological
1323 monthly mean value for January. This helps answer the question "is my model a better source of information than
1324 using the average past conditions from this time of year?".

1325 **Connectome.** The connections between nodes in a neural network. Examples include fully-connected, partially-
1326 connected, skip-layer connections, recurrent connections and others. The 'wiring diagram' for the network.

1327 **Convolutional neural network.** A neural network architecture commonly applied to images which utilises a
1328 convolutional (spatially connected) kernel applied in a sliding window fashion with a narrow receptive field to
1329 encourage the network to generalise from fine scale structure to higher levels of abstraction.

1330 **Data augmentation.** The practice of modifying input data in supervised learning to produce additional examples.
1331 This can make networks more robust to new inputs and address issues of brittleness to adversarial examples. An
1332 example of data augmentation is using rotated or reflected versions of the same image as independent training samples.

1333 **Data driven.** A generalised term used to indicate a primary reliance or dependence on the collection or analysis of
1334 data. Used in contrast to process driven or theory driven.

1335 **Decision tree.** A tree-like, or flowchart-like, branching model representing a series of decisions and their possible
1336 consequences. Each internal node represents a 'test' (i.e. decision threshold) and each leaf node represents a class label
1337 or collection of possible outcomes.

1338 **Deep NN.** A neural network with many layers. Deeper, thinner networks have generally been more popular in recent
1339 times than wider, shallower ones but this is not always the case (see e.g. Zagoruyko & Komodakis, 2016)

1340 **DeepONet.** A neural network architecture relying on universal approximation theorem to train a neural network to
1341 represent a mathematical operation (the operator), such as a partial differential equation or dynamic system.

1342 **Discriminator model.** A model which distinguishes or discriminates between synthetic data and real-world
1343 observations. Often used in conjunction with a generator. In this case, the overall goal is to produce a generator which
1344 is capable of fooling the discriminator, producing highly realistic images. This process is used in Generative
1345 Adversarial Networks.

1346 **Dropout layer.** A neural network layer which is only partially connected, often with a stochastic dropout chance. This
1347 has been shown experimentally to improve neural network robustness in many architectures by reducing overfitting.

1348 **Epoch.** A single complete training pass through all available training data, e.g. learning from all samples, or learning
1349 from all mini-batches, according to the training strategy. Multiple training epochs will typically be utilised although
1350 alternative strategies do exist.

1351 **Feed-forward network.** A neural network composed of distinct 'layers', where the outputs of one layer never feed
1352 back into earlier layers. This avoids the needs for any iterative solver approaches and results in a very computationally
1353 efficient 'forward pass'.

1354 **Generative adversarial network.** A two-part neural network architecture comprising a generator and a discriminator,
1355 which are co-trained to produce realistic outputs which are hard to distinguish from real-world data. The discriminator
1356 replaces the traditional loss function.

1357 **Generator model.** A model which produces a synthetic example of a particular class, such as a synthetic image or
1358 synthetic language. Examples include language or image generation. These are used as part of Generative Adversarial
1359 Networks among other applications.

1360 **Global receptive field.** Where every part of the input region can influence or stimulate a response in a model (e.g. a
1361 fully-connected neural network).

1362 **GPU.** Graphical Processing Unit. A hardware device specialised for fast matrix operations, originally created to
1363 support computer graphics, particularly for games.

1364 **Gradient boosted decision tree.** Also referred to as extreme gradient boosting. A random forest architecture which
1365 combines gradient boosting with decision tree ensembles.

1366 **Gradient boosting.** An approach to model training where each additional ensemble member attempts to predict the
1367 cumulative errors of previously trained members.

1368 **Graph neural network.** A class of neural networks designed to process data which is described by a graph (or
1369 tree/network) data structure. See Scarselli et al. (2008), Kipf & Welling (2016), and Battaglia et al. (2018) for more
1370 information and examples.

1371 **Hidden layer.** A layer which is intermediate between the input layer and the output layer of a network or tree structure.
1372 Hidden layers may be used to encode 'hidden variables' which are latent to a problem but not able to be directly
1373 observed.

1374 **Hierarchical temporal aggregation.** A mechanism of composing neural networks which are trained for different lead
1375 times to produce an optimal prediction at all time horizons.

1376 **Hierarchical temporal memory.** Fundamentally different to hierarchical temporal aggregation. A complex deep
1377 learning architecture which uses time-adjacency pooling.

1378 **Hyperparameter.** A parameter which is not derived via training. Examples include the learning rate and the model
1379 topology.

1380 **Hyperparameter search (or Hyperparameter optimization).** The process of determining optimal hyperparameters.
1381 This term may also be used to encompass the model selection problem. This process is automated in some cases.

1382 **Input layer.** A layer which is composed of input nodes. Typically machine learning models will have one input layer
1383 at depth zero (i.e. with no preceding layers) and no input nodes at greater depths.

1384 **Input node.** A node which represents an input or observed value.

1385 **K-fold cross-validation.** A process of changing the validation and test data partitions during different iterations of
1386 training. This allows more of the training and validation data to be used while minimising overfitting. Some definitions
1387 include test data in this process but that is not ideal as the final test is no longer statistically independent.

1388 **Keras.** A streamlined API for creating neural networks, integrated with Tensorflow. Originally built on the Theano
1389 framework for general mathematical evaluation. PyTensor and Aesara are related packages.

1390 **Kernel trick.** For data sets which are not linearly separable, first multiplying the data by a nonlinear function in a
1391 higher dimension can result in a linearly separable higher-dimensional data set to which a simpler method can be used
1392 to model the data.

1393 **Knowledge based systems.** A broad term from artificial intelligence meaning a system which that uses reasoning and
1394 a knowledge base to support decision making. Knowledge is represented explicitly and a reasoning or inference engine
1395 is used to arrive at new knowledge.

1396 **Layer.** In tree or feed-forward network structures (e.g. decision trees and feed-forward neural networks), a layer refers
1397 to the set of nodes at the same depth within a network.

1398 **Leaf node.** Aka output node. A node which does not have any child nodes.

1399 **Long short term memory network.** A recurrent neural network architecture which processes sequences of tokens
1400 utilising a 'memory' component which can store information from tokens early in a sequence for use in prediction of
1401 tokens much later in a sequence. Typical applications include language prediction and time-series prediction of many
1402 kinds.

1403 **Loss function** (also known as target function, training function, objective function, penalty score, error function,
1404 heuristic function, minimisation function). A differentiable function which is well-behaved, such that smaller values
1405 represent better model performance and larger values represent worse performance. An example would be the root-
1406 mean-squared-error of a prediction compared to the truth or target value.

1407 **Mini batch.** A subset or 'mini batch' of the training data. Utilised for multiple reasons, including computational
1408 efficiency and to reduce overfitting. Aggregate error over a mini-batch is be learned rather than per-sample errors.
1409 This is the typical contemporary approach. See also training batch for in-depth discussion.

1410 **Neural network.** A composition of 'input nodes', 'connections', 'nodes', 'layers', 'output layers' and 'activation
1411 functions' which are capable of complex modelling tasks. Originally designed to simulate human neural functioning
1412 and subsequently applied to a range of applications.

1413 **Node. Aka vertex.** A small data structure in a network, tree or graph structure which is connected by edges. A node
1414 may represent a real-world value (such as a location) or an abstract value (such as in a neural network), or a decision
1415 threshold (such as in a decision tree).

1416 **Normalisation.** A technique applied in many areas of mathematics, science and statistics which is also very important
1417 to machine learning and neural networks. In a general sense, this refers to expressing values within a standard range.
1418 Very often, the range of expected values is mapped onto the range 0 to 1, to allow physical variables with different
1419 measurement units to be compared on equal scale. Such normalisation may be linear or nonlinear, according to a
1420 simple or more complex function, and either drawn from known physical limits or from the variation observed in the
1421 data itself.

1422 **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
1423 machine learning for language processing to create a word or feature embedding in a process called tokenisation or
1424 encoding. The length of the vector is commonly equal to the number of categories or symbols.

1425 **Output layer.** A layer which comprises the leaf nodes or output nodes of a tree or network.

1426 **Perceptron.** A single-layer neural network architecture for supervised learning of binary classification. Originally
1427 built as an electronic hardware device encoding weights with potentiometers and learning with motors. A multi-layer
1428 perceptron is the same thing as an ordinary neural network.

1429 **Persistence.** Refers to the practice of treating some past observation or reanalysis (usually immediately prior to the
1430 starting point of the prediction period) as the future prediction and "persisting" this one state forward to every
1431 prediction lead time. The predictive model is then compared to this persistence prediction, essentially assessing the
1432 performance of the model against a steady state prediction. This, along with climatology, is often used as a baseline
1433 or bare minimum prediction to beat (i.e., a prediction better than persistence could be considered skilful vs
1434 persistence). This answers the question "is my model a better source of information than using what happened just
1435 before now?".

1436 **Physically-informed machine learning. Also known as physics-informed machine learning.** Machine learning is
1437 considered physically informed when some aspect of physics is included in any way. Examples include adding a
1438 physical component to the loss function (e.g. to enforce conservation of physical properties) or using an activation
1439 function with physically realistic properties.

1440 **Predictive step, forward pass, evaluation.** The process of calculating a model prediction from a set of input
1441 conditions. Distinct from the training phase or back-propagation step.

1442 **PyTorch.** A widely adopted framework for neural networks in Python.

1443 **Random forest.** An architecture based on decision tree ensembles where each decision tree is initialised semi-
1444 randomly and an average of all models is used for prediction. This is typically more accurate than a single decision
1445 tree but less accurate than a gradient-boosted decision tree and so is now less-used. The term random forest is still
1446 commonly used when in fact the implementation is a gradient boosted decision tree.

1447 **Receptive field.** The size or extent of a region in the input which can influence or stimulate a response in a model,
1448 e.g. the size of a convolutional kernel, the size of a sliding window

1449 **Rectified Linear Unit (ReLU).** An activation function commonly used in DNNs. Defined as $\max(0, X)$. This function
1450 is used as it is computationally cheap and avoids problems of vanishing gradients.

1451 **Recurrent network.** A neural network which does pass the output from nodes of the network back into the input of
1452 others. Infinite recurrence is avoided by setting a specific number of iterations for the recurrence. These are often
1453 depicted in diagrams as separate layers but the implementation is through internal recurrent connections.

1454 **Regression.** A model which attempts to diagnose or predict an exact value by statistically relating example input
1455 values to desired values.

1456 **Relevance vector machine.** A sparse Bayesian model utilising the kernel trick in similar fashion to a support vector
1457 machine.

1458 **Representation error.** Error which is introduced due to the inexactness of representing the real world in the model
1459 belief state. Examples may include topography smoothing, point-to-grid translations, model grid distortions near the
1460 poles, or the exclusion of physical characteristics which are not primary to the model.

1461 **Residual neural network (ResNet).** A very influential and innovative convolutional NN architecture which uses a
1462 similar concept to gradient boosting. Each layer of the deep network is taken to predict the residual error from the
1463 previous layers, with skip-connections from earlier layers allowing the training to occur without the issue of vanishing
1464 gradients.

1465 **Sample.** A single training example (e.g. a row of data).

1466 **Scale invariance.** A feature of a system, problem or model which means the results and behaviour are the same at any
1467 scale (e.g., the behaviour does not change if the inputs are multiplied by a common factor).

1468 **Scikit-learn.** A popular Python library for machine learning which extends the SciPy framework.

1469 **Sharding.** Refers to dividing the training of a neural network across multiple GPUs or nodes. This can be done using
1470 data sharding, whereby each GPU or node trains on a subset of the data to allow training parallelism, or model sharding
1471 where a single model is partitioned across multiple GPUs to allow a larger neural network than could be allocated in
1472 memory on a single GPU. One example could be assigning a small number neural network layers to each GPU which
1473 could then work in sequence to operate on a very large network.

1474 **(Stochastic) Gradient descent.** An algorithm by which a neural network is trained using increasingly fine-scale
1475 adjustments to optimise the accuracy of network prediction. Utilised to find the local minimum of a differentiable
1476 function.

1477 **Supervised learning.** Machine learning is considered 'supervised' when the data is labelled according to a category
1478 or target value. Classification data have an explicit labelled category. Regression data have an explicit value which is
1479 being predicted for.

1480 **Support vector machine.** A classification model based on finding a hyperplane to separate data utilising the kernel
1481 trick.

1482 **Tensor.** Can be considered as a dense multi-dimensional array or matrix.

1483 **Tensorflow.** A widely adopted framework for neural networks in Python.

1484 **Test/train/validate split.** Available data is split into three portions. The training data is evaluated and used to update
1485 model weights. Validation data is evaluated during training and may be used for hyper-parameter search or to guide
1486 the researcher. Test data is independent (typically well-curated) data used for gold standard evaluation. In reality,
1487 validation data is sometimes used as test data, but this is not good practice. There are many considerations for
1488 test/train/validate splitting, such as statistical independence, representation of all classes, and bias. It is important to
1489 consider what the model is generalising "from" and "to", and ensuring appropriate examples are present in the training
1490 data and appropriate examples are reserved for validation and test.

1491 **Token.** Tokenisation the process of mapping a symbolic or categorical sequence to a numerical representation which
1492 is suited to a sequence-based machine learning model. Commonly, a vector representation will be utilised for the token
1493 form. In language processing, either characters or words may be represented as tokens depending on the approach.

1494 **Top Hat function.** A filter or function which has a rectangular shape resembling the cross-section of a top hat. One
1495 of the simplest functions used for convolutional operations, it can be defined as one constant value in a given bounded
1496 range, and another smaller constant value outside that range.

1497 **TPU.** Tensor Processing Unit. A hardware device specialised for artificial intelligence and machine learning
1498 applications, in particular neural network operations.

1499 **Training batch (or simply batch).** Multiple definitions apply and the use the term has evolved over time. Originally
1500 used in the context of learning from offline or saved historical data as opposed to online or realtime novel data. In this
1501 definition, the training batch is the saved data and refers to the whole training set. For example, a robot exploring a
1502 new environment in real-time must use an online learning technique and could not utilise batch training to map the
1503 unseen terrain. In more recent use, particularly in the areas of neural network learning, the offline saved data may be
1504 split into one or more batches (subsets). If one batch (the batch is the entire training set) is used, the aggregate errors
1505 for the entire training set are used to update the model weights and biases, and the learning algorithm is called batch
1506 gradient descent. If each example is presented individually, this is called online training (even when historical saved
1507 data is being used), the weights and biases are updated for from each individual example, and the algorithm used is
1508 stochastic gradient descent. If the data is divided into multiple batches, this is often referred to equivalently as mini
1509 batches. The weights and biases are aggregated over each mini batch. This is the most common contemporary
1510 approach, as it reduces overfitting and is a good balance of training accuracy, avoiding local minima, and
1511 computational efficiency.

1512 **Transfer learning.** The process of training a model first on a related problem, and then conducting further training
1513 on a more specific problem. Examples could be training a model first in one geographical region and then in another;
1514 or training first at a low resolution then subsequently at a high resolution. This is frequently done to reduce training
1515 computation cost for similar problems by re-using the trained weights from a well-performing source model, or to
1516 overcome a problem of limited data availability by using multiple data sources.

1517 **Transformer network.** A token-sequence architecture which is capable of handling long-range dependencies.
1518 Initially applied to language processing, it has found effective application in image processing as an alternative to
1519 convolutional architectures.

1520 **Translation invariance.** A feature of a system, problem or model which means the results and behaviour are the same
1521 after any spatial translation (i.e., the behaviour does not change if the inputs are shifted spatially to a new location).
1522 **U-Net.** A type of convolutional neural network developed for biomedical image segmentation which has found broad
1523 application. In the contracting part of the network spatial information is reduced while feature information is increased.
1524 In the expanding part of the network, feature information is used to inform high-resolution segmentation. The name
1525 derives from the diagrammatic shape of the network forming a "U".
1526 **Unsupervised learning.** Machine learning is considered 'unsupervised' when data is unlabelled. Examples include
1527 clustering, association and dimensionality reduction.
1528 **Vanishing Gradient.** At the extremes, nonlinear functions used to calculate gradients can result in gradient values
1529 which are effectively zero. These small or zero values, once present in the weights and biases of a neural network, can
1530 entirely suppress information which would in fact be useful, and result in a local minima from which training cannot
1531 recover. This is particularly relevant to long token-series when long-distance connections are relevant. A variety of
1532 techniques including alternative activation functions, training weight decay, skip connections and attention
1533 mechanisms may each or all be utilised to ameliorate this issue.
1534 **Weights and biases.** The parameter values for each neuron which represent the weighting factors to apply to the input
1535 values, plus an overall bias value for the node.
1536 **XGBoost.** A popular Python library for gradient boosted decision trees.

1537
1538 **Appendix A: Table Summary of Model Architectures cited in this paper.**
1539 This table includes all references from this review except for: seminal ML papers that are on new ML methods (e.g., foundational
1540 ML papers), review papers, any paper cited that concerns a topic which is out of scope (e.g., nowcasting), and any other paper
1541 which does not present a new method directly applicable to weather and climate modelling.

Author(s)	Year	Category	Approach
Ackmann et al	2020	Fully connected NN	Preconditioner
Alemohammad et al	2017	Fully connected NN	Variable estimation
Andersson et al	2021	Convolutional NN	Prediction
Arcomano et al	2022	Reservoir computing	Alongside-model bias corrector
Atkinson	2020	Baysean type NN	PDE solver
Bar-Sinai	2019	Convolutional NN	PDE solver
Battaglia et al	2018	Graph NN	Method paper
Beucler et al	2019	Physics Informed NN	Convective paramterisation
Beucler et al	2021	Physics Informed NN	Convective paramterisation
Bhattacharya et al	2021	Fully connected NN	PDE solver
Bi et al	2022	Mixed/Custom NN	Pure ML atmospheric model
Bihlo & Popovych	2022	Physics Informed NN	PDE solver
Bolton and Zanna	2019	Convolutional NN	Parametrization
Brenowitz & Bretherton	2018	Fully connected NN	Parametrization

Brenowitz & Bretherton	2019	Fully connected NN	Parametrization
Brenowitz et al.	2020	Fully connected NN	Parametrization
Brenowitz et al	2020	Decision tree-based, Fully connected NN	ML model intercomaprison
Brenowitz et al	2022	Recurrent NN	Parametrization
Chaney et al	2016	Decision tree-based	Interpolation
Chantry et al	2021	Fully connected NN	Parametrization
Chattopadhyay et al	2020	Fully connected NN, Recurrent NN	Super parametrization
Chevallier et al	1998	Fully connected NN	Parametrization
Chi & Kim	2017	Fully connected NN, Recurrent NN	Prediction
Clare et al	2021	ResNet	Emulation (probabilistic)
Dagon et al	2020	Fully connected NN	Emulation
de Bézenac et al	2017	GAN	Prediction, model evaluation
Deuben and Bauer	2018	Fully connected NN	Replacement
Flora et al	2022	Decision tree-based, Logistic regression	Assessment of explainability techniques
Fuhg et al	2022	Physics Informed NN	PDE solver
Gagne et al	2019	Decision tree-based	Parametrization
Gagne et al	2020	GAN	Parametrization (probabilistic)
Gagne et al	2020	GAN, Fully connected NN	Parametrization
George et al	2008	Mixed/Custom non-NN	Preconditioner
Gettelman et al	2021	Fully connected NN	Emulation
Ham et al	2019	Convolutional NN	Prediction
Ham et al	2021	Convolutional NN	Prediction
Han et al	2020	ResNet	Parametrization
Harder et al	2022	Fully connected NN	Emulation
He et al	2022	Decision tree-based	Parametrization
Holloway & Chen	2007	Fully connected NN	Preconditioner and PDE solver selection
Horvat & Roach	2022	Fully connected NN	Parametrization
Hu et al	2023	Mixed/Custom NN	Pure ML atmospheric model
Huang et al	2016	SVM	Preconditioner
Kapp-Schwoerer et al	2020	Convolutional NN	Semantic segmentation
Karunasinghe & Liong	2006	Fully connected NN	Chaotic timeseries prediction
Keisler	2022	Graph NN	Replacement
Kim et al	2022	Mixed/Custom NN	Prediction
Kochkov et al	2021	Convolutional NN	PDE solver
Krasnopolsky et al	2002	Fully connected NN	Emulation
Krasnopolsky et al	2005	Fully connected NN	Emulation
Krasnopolsky	2013	Fully connected NN	Parametrization (probabilistic)

Kuefler & Chen	2008	Mixed/Custom non-NN	Linear system solver
Ladický et al	2015	Decision tree-based	PDE solver
Lam et al	2022	Mixed/Custom NN	Pure ML atmospheric model
Lanthaler et al	2022	Neural Operator	PDE solver
Leufen & Schadler	2019	Fully connected NN	Parameterization
Li et al	2020	Graph NN	PDE solver
Li et al	2020	Neural Operator	PDE solver
Li et al	2020	Neural Operator	PDE solver
Lopez-Gomez et al	2023	Convolutional NN	Prediction
Lu et al	2020	Neural Operator	PDE solver
Meyer et al	2022	Fully connected NN	Emulation
Moishin et al	2021	Convolutional Recurrent NN	Prediction
Mooers et al	2021	Fully connected NN	Emulation
Mudigonda et al	2017	Mixed/Custom NN	Object detection
Nelsen & Stuart	2021	Random Feature Model	PDE solver
Nguyen et al	2023	Mixed/Custom NN	Pure ML atmospheric model
O'Brien et al	2020	Bayesian model	Object detection
O'Gorman & Dwyer	2018	Decision tree-based	Emulation
O'Leary et al	2022	Fully connected NN	PDE solver
Ott et al	2020	Fully connected NN	Emulation
Pan et al	2020	Decision tree-based	Parameterisation
Patel et al	2021	Neural Operator	PDE solver
Pathak et al	2022	Mixed/Custom NN	Pure ML atmospheric model
Peairs & Chen	2011	Mixed/Custom non-NN	PDE solver
Pelissier et al	2020	Mixed/Custom non-NN	Hybrid model corrector
Prabhat et al	2021	Convolutional NN	Object detection
Psaros et al	2023	Neural Operator, Physics Informed NN	PDE solver
Rasp	2020	Fully connected NN	Emulation
Rasp et al	2018	Fully connected NN	Emulation
Rasp et al	2020	Fully connected NN, Linear regression	Pure ML atmospheric model
Rasp & Thuerey	2021	ResNet	Pure ML atmospheric model
Rizzuti et al	2019	Convolutional NN	NN based corrector step in PDE solver
Rosier et al	2023	Mixed/Custom NN	Prediction
Ross et al.	2023	Genetic programming, Linear regression, Convolutional NN	Intercomparison of methods to learn parameterisations from data
Rupe et al	2023	Mixed/Custom non-NN	Object detection
Sawada	2020	Regression	Emulation
Scher	2018	Convolutional NN	Emulation

Scher and Messori	2019	Convolutional NN	Emulation
Taylor & Feng	2022	Convolutional NN	Prediction
Tompson et al	2017	Convolutional NN	PDE solver
Toms et al	2020	Fully connected NN	NN interpretability
Ukkonen & Mäkelä	2019	Decision tree-based, Logistic Regression, Fully connected NN	Paramterisation
Ukkonen et al	2020	Fully connected NN	Emulation
Vlachas et al	2018	Recurrent NN	Pure ML baseline model
Wang et al	2021	Neural Operator	PDE solver
Wang et al	2022	ResNet	Parametrization
Wang et al	2022	Physics Informed NN	PDE solver
Watt-Meyer et al	2021	Decision tree-based	Nudging
Watson-Parris et al	2022	Gaussian Process, Decision tree-based, Mixed/Custom NN	Pure ML baseline model
Weyn et al	2019	Convolutional NN	Pure ML atmospheric model
Weyn et al	2020	Convolutional NN	Pure ML atmospheric model
Weyn et al	2021	Convolutional NN	Pure ML atmospheric model
Wikner et al	2020	Reservoir computing	Alongside-model bias corrector
Wu & Xiu	2020	ResNet	Learning PDE operators
Yamada et al	2018	Convolutional NN	Preconditioner
Yang et al	2016	Fully connected NN	PDE solver
Yeo et al	2021	Recurrent NN	Dynamical system simulation
Yuval & O’Gorman	2020	Decision tree-based	Emulation
Yuval et al	2021	Fully connected NN	Emulation
Zanna and Bolton	2020	Convolutional NN, Relevance vector machine	Parametrization and equation discovery
Zhao et al	2019	Fully connected NN	Paramterisation
Zhao et al	2019	Physics Informed NN	Paramterisation
Zhong et al	2023	Fully connected NN, Recurrent NN	Emulation

1542

1543

1544 **Code Availability**

1545 No code was used in the preparation of this review.

1546 **Data Availability**

1547 No data was processed in the preparation of this review except for the list of ML model types by cited paper, which
1548 is provided in the appendix.

1549 **Author Contribution**

1550 COdBD researched and wrote Sections 3, 4, 5, 6 and 7, and provided review of sections 8, 10, and the glossary. TL
1551 researched and wrote sections 8, 10, and the glossary, and provided review of sections 3, 4, 5, 6, and 7. COdBD and
1552 TL researched and co-wrote sections 1, 2, 9, 11, 12, and the Appendix.

1553 **Competing Interests**

1554 The authors declare that they have no conflict of interest.

1555 **Acknowledgements**

1556 The authors would like to thank Bethan White, Harrison Cook, Tom Dunstan and Karina Williams for their very
1557 helpful reviews of early versions of this manuscript. We also would like to wholeheartedly thank the referees for their
1558 extremely helpful, positive and well considered feedback and suggestions. Their input has greatly improved this
1559 review. Finally, we would like to acknowledge and thank the people who contacted us with comments, suggestions,
1560 and advice on the preprint versions of this review. All of the input was valuable, and greatly appreciated.

1561 **References**

- 1562 Ackmann, J., Düben, P. D., Palmer, T. N., & Smolarkiewicz, P. K. Machine-learned preconditioners for linear solvers
1563 in geophysical fluid flows. *arXiv preprint arXiv:2010.02866*. <https://doi.org/10.48550/arXiv.2010.02866>. 6
1564 October 2020.
- 1565 Adams, S. V., Ford, R. W., Hambley, M., Hobson, J. M., Kavčič, I., Maynard, C. M., ... & Wong, R. LFRic: Meeting
1566 the challenges of scalability and performance portability in Weather and Climate models. *J PARALLEL*
1567 *DISTR COM*, 132, 383-396. <https://doi.org/10.1016/j.jpdc.2019.02.007>. 2019.
- 1568 Alemohammad, S. H., Fang, B., Konings, A. G., Aires, F., Green, J. K., Kolassa, J., ... & Gentine, P. Water, Energy,
1569 and Carbon with Artificial Neural Networks (WECANN): a statistically based estimate of global surface
1570 turbulent fluxes and gross primary productivity using solar-induced fluorescence. *BIOGEOSCIENCES*,
1571 14(18), 4101-4124. <https://doi.org/10.5194/bg-14-4101-2017>. 2017.
- 1572 Andersson, T. R., Hosking, J. S., Pérez-Ortiz, M., Paige, B., Elliott, A., Russell, C., ... & Shuckburgh, E. Seasonal
1573 Arctic sea ice forecasting with probabilistic deep learning. *NAT COMMUN*, 12(1), 5124.
1574 <https://doi.org/10.1038/s41467-021-25257-4>. 2021.
- 1575 Arcomano, T., Szunyogh, I., Wikner, A., Pathak, J., Hunt, B. R., & Ott, E. A Hybrid Approach to Atmospheric
1576 Modeling That Combines Machine Learning With a Physics-Based Numerical Model. *J ADV MODEL*
1577 *EARTH SY*, 14(3), e2021MS002712. <https://doi.org/10.1029/2021MS002712>. 2022.
- 1578 Atkinson, S. Bayesian hidden physics models: Uncertainty quantification for discovery of nonlinear partial differential
1579 operators from data. *arXiv preprint arXiv:2006.04228*. <https://doi.org/10.48550/arXiv.2006.04228>. 7 June
1580 2020.

1581 Bar-Sinai, Y., Hoyer, S., Hickey, J., & Brenner, M. P. Learning data-driven discretizations for partial differential
1582 equations. *Proceedings of the National Academy of Sciences*, 116(31), 15344-15349.
1583 <https://doi.org/10.1073/pnas.1814058116>. 2019.

1584 Battaglia, P. W., Hamrick, J. B., Bapst, V., Sanchez-Gonzalez, A., Zambaldi, V., Malinowski, M., ... & Pascanu, R.
1585 Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*.
1586 <https://doi.org/10.48550/arXiv.1806.01261>. 4 June 2018.

1587 Beucler, T., Rasp, S., Pritchard, M., & Gentine, P. Achieving conservation of energy in neural network emulators for
1588 climate modeling. *arXiv preprint arXiv:1906.06622*. <https://doi.org/10.48550/arXiv.1906.06622>. 15 June
1589 2019.

1590 Beucler, T., Pritchard, M., Rasp, S., Ott, J., Baldi, P., & Gentine, P. Enforcing analytic constraints in neural
1591 networks emulating physical systems. *PHYS REV LETT*, 126(9), 098302.
1592 <https://doi.org/10.1103/PhysRevLett.126.098302>. 2021.

1593 Bhattacharya, K., Hosseini, B., Kovachki, N. B., & Stuart, A. M. Model reduction and neural networks for parametric
1594 PDEs. *arXiv preprint arXiv:2005.03180*. <https://doi.org/10.48550/arXiv.2005.03180>. 7 May 2020.

1595 Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X., & Tian, Q. Pangu-Weather: A 3D High-Resolution Model for Fast and
1596 Accurate Global Weather Forecast. *arXiv preprint arXiv:2211.02556*.
1597 <https://doi.org/10.48550/arXiv.2211.02556>. 3 November 2022.

1598 Bihlo, A., & Popovych, R. O. Physics-informed neural networks for the shallow-water equations on the sphere. *J*
1599 *COMPUT PHYS*, 456, 111024. <https://doi.org/10.1016/j.jcp.2022.111024>. 2022.

1600 Bolton, T., & Zanna, L. Applications of deep learning to ocean data inference and subgrid parameterization. *J ADV*
1601 *MODEL EARTH SY*, 11(1), 376-399. <https://doi.org/10.1029/2018MS001472>. 2019.

1602 Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32.

1603 Brenowitz, N. D., & Bretherton, C. S. Prognostic validation of a neural network unified physics parameterization.
1604 *GEOPHYS RES LETT*, 45(12), 6289-6298. <https://doi.org/10.1029/2018GL078510>. 2018.

1605 Brenowitz, N. D., & Bretherton, C. S. Spatially extended tests of a neural network parametrization trained by coarse-
1606 graining. *J ADV MODEL EARTH SY*, 11(8), 2728-2744. <https://doi.org/10.1029/2019MS001711>. 2019.

1607 Brenowitz, N. D., Beucler, T., Pritchard, M., & Bretherton, C. S. Interpreting and stabilizing machine-learning
1608 parametrizations of convection. *J ATMOS SCI*, 77(12), 4357-4375. <https://doi.org/10.1175/JAS-D-20-0082.1>. 2020.

1610 Brenowitz, N. D., Henn, B., McGibbon, J., Clark, S. K., Kwa, A., Perkins, W. A., ... & Bretherton, C. S. Machine
1611 learning climate model dynamics: Offline versus online performance. *arXiv preprint arXiv:2011.03081*.
1612 <https://doi.org/10.48550/arXiv.2011.03081>. 5 November 2020.

1613 Brenowitz, N. D., Perkins, W. A., Nugent, J. M., Watt-Meyer, O., Clark, S. K., Kwa, A., ... & Bretherton, C. S.
1614 Emulating Fast Processes in Climate Models. *arXiv preprint arXiv:2211.10774*.
1615 <https://doi.org/10.48550/arXiv.2211.10774>. 19 November 2022.

1616 Carranza-García, M., García-Gutiérrez, J., & Riquelme, J. C. A framework for evaluating land use and land cover
1617 classification using convolutional neural networks. *REMOTE SENS-BASEL*, 11(3), 274.
1618 <https://doi.org/10.3390/rs11030274>. 2019.

1619 Chaney, N. W., Herman, J. D., Ek, M. B., & Wood, E. F. Deriving global parameter estimates for the Noah land
1620 surface model using FLUXNET and machine learning. *J GEOPHYS RES-ATMOS*, 121(22), 13-218.
1621 <https://doi.org/10.1002/2016JD024821>. 2016.

1622 Chantry, M., Christensen, H., Dueben, P., & Palmer, T. Opportunities and challenges for machine learning in weather
1623 and climate modelling: hard, medium and soft AI. *Philosophical Transactions of the Royal Society A*,
1624 379(2194), 20200083. <https://doi.org/10.1098/rsta.2020.0083>. 2021.

1625 Chantry, M., Hatfield, S., Dueben, P., Polichtchouk, I., & Palmer, T. Machine learning emulation of gravity wave
1626 drag in numerical weather forecasting. *J ADV MODEL EARTH SY*, 13(7), e2021MS002477.
1627 <https://doi.org/10.1029/2021MS002477>. 2021.

1628 Chase, R. J., Harrison, D. R., Burke, A., Lackmann, G. M., & McGovern, A. (2022a). A Machine Learning Tutorial
1629 for Operational Meteorology, Part I: Traditional Machine Learning. *arXiv preprint arXiv:2204.07492*.
1630 <https://doi.org/10.48550/arXiv.2204.07492>. 15 April 2022

1631 Chase, R. J., Harrison, D. R., Lackmann, G., & McGovern, A. A Machine Learning Tutorial for Operational
1632 Meteorology, Part II: Neural Networks and Deep Learning. *arXiv preprint arXiv:2211.00147*.
1633 <https://doi.org/10.48550/arXiv.2211.00147>. 31 October 2022.

1634 Chattopadhyay, A., Subel, A., & Hassanzadeh, P. Data-driven super-parameterization using deep learning:
1635 Experimentation with multiscale Lorenz 96 systems and transfer learning. *J ADV MODEL EARTH SY*,
1636 12(11), e2020MS002084. <https://doi.org/10.1029/2020MS002084>. 2020.

1637 Chevallier, F., Chéruy, F., Scott, N. A., & Chédin, A. A neural network approach for a fast and accurate computation
1638 of a longwave radiative budget. *J APPL METEOROL*, 37(11), 1385-1397. [https://doi.org/10.1175/1520-0450\(1998\)037%3C1385:ANNAFA%3E2.0.CO;2](https://doi.org/10.1175/1520-0450(1998)037%3C1385:ANNAFA%3E2.0.CO;2). 1998.

1640 Chi, J., & Kim, H. C. Prediction of arctic sea ice concentration using a fully data driven deep neural network. *REMOTE*
1641 *SENS-BASEL*, 9(12), 1305. <https://doi.org/10.3390/rs9121305>. 2017.

1642 Clare, M. C., Jamil, O., & Morcrette, C. J. (2021). Combining distribution-based neural networks to predict weather
1643 forecast probabilities. *Q J ROY METEOR SOC*, 147(741), 4337-4357. <https://doi.org/10.1002/qj.4180>.
1644 2021.

1645 Conley, A. J., Garcia, R., Kinnison, D., Lamarque, J. F., Marsh, D., Mills, M., ... & Taylor, M. A. Description of the
1646 NCAR community atmosphere model (CAM 5.0). NCAR technical note, 3. 2012.

1647 Cuomo, S., Di Cola, V. S., Giampaolo, F., Rozza, G., Raissi, M., & Piccialli, F. Scientific Machine Learning through
1648 Physics-Informed Neural Networks: Where we are and What's next. *arXiv preprint arXiv:2201.05624*.
1649 <https://doi.org/10.48550/arXiv.2201.05624>. 14 January 2022.

1650 Dagon, K., Sanderson, B. M., Fisher, R. A., & Lawrence, D. M. A machine learning approach to emulation and
1651 biophysical parameter estimation with the Community Land Model, version 5. *Advances in Statistical*

1652 Climatology, Meteorology and Oceanography, 6(2), 223-244. <https://doi.org/10.5194/ascmo-6-223-2020>.
1653 2020.

1654 De Bézenac, E., Pajot, A., & Gallinari, P. *Towards a hybrid approach to physical process modeling*. Technical report.
1655 2017.

1656 de Witt, C. S., Tong, C., Zantedeschi, V., De Martini, D., Kalaitzis, F., Chantry, M., ... & Bilinski, P. RainBench:
1657 towards global precipitation forecasting from satellite imagery. arXiv preprint arXiv:2012.09670.
1658 <https://doi.org/10.48550/arXiv.2012.09670>. 17 December 2020.

1659 Deng, J., Dong, W., Socher, R., Li, L. J., Li, K., & Fei-Fei, L. Imagenet: A large-scale hierarchical image database.
1660 PROC CVPR IEEE (pp. 248-255). Ieee. <https://doi.org/10.1109/CVPR.2009.5206848>. 2009.

1661 Digra, M., Dhir, R., & Sharma, N. Land use land cover classification of remote sensing images based on the deep
1662 learning approaches: a statistical analysis and review. ARAB J GEOSCI, 15(10), 1003.
1663 <https://doi.org/10.1007/s12517-022-10246-8>. 2022.

1664 Dijkstra, H. A., Petersik, P., Hernández-García, E., & López, C. The application of machine learning techniques to
1665 improve El Niño prediction skill. AIP CONF PROC, 153. <https://doi.org/10.3389/fphy.2019.00153>. 2019.

1666 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... & Houlsby, N. An image is
1667 worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929.
1668 <https://doi.org/10.48550/arXiv.2010.11929>. 22 October 2020.

1669 Dueben, P. D., & Bauer, P. Challenges and design choices for global weather and climate models based on machine
1670 learning. GEOSCI MODEL DEV, 11(10), 3999-4009. <https://doi.org/10.5194/gmd-11-3999-2018>. 2018.

1671 Dueben, P. D., Schultz, M. G., Chantry, M., Gagne, D. J., Hall, D. M., & McGovern, A. Challenges and Benchmark
1672 Datasets for Machine Learning in the Atmospheric Sciences: Definition, Status, and Outlook. *Artificial*
1673 *Intelligence for the Earth Systems*, 1(3), e210002. <https://doi.org/10.1175/AIES-D-21-0002.1>. 2022.

1674 ECMWF. (2018). IFS documentation (cy45r1). Retrieved from <https://www.ecmwf.int/en/publications/ifs->
1675 [documentation](https://www.ecmwf.int/en/publications/ifs-documentation) accessed 7th February 2023

1676 Eyring, V., Bony, S., Meehl, G. A., Senior, C. A., Stevens, B., Stouffer, R. J., & Taylor, K. E. Overview of the Coupled
1677 Model Intercomparison Project Phase 6 (CMIP6) experimental design and organization. GEOSCI MODEL
1678 DEV, 9(5), 1937-1958. <https://doi.org/10.5194/gmd-9-1937-2016>. 2016.

1679 Flora, M., Potvin, C., McGovern, A., & Handler, S. Comparing Explanation Methods for Traditional Machine
1680 Learning Models Part 2: Quantifying Model Explainability Faithfulness and Improvements with
1681 Dimensionality Reduction. arXiv preprint arXiv:2211.10378. <https://doi.org/10.48550/arXiv.2211.10378>. 18
1682 November 2022.

1683 Friedman, J. H. Greedy function approximation: a gradient boosting machine. ANN STAT, 1189-1232. 2001.

1684 Fuhg, J. N., Kalogeris, I., Fau, A., & Bouklas, N. Interval and fuzzy physics-informed neural networks for uncertain
1685 fields. PROBABILIST ENG MECH, 68, 103240. 2022.

1686 Gagne, D. J., McCandless, T., Kosovic, B., DeCastro, A., Loft, R., Haupt, S. E., & Yang, B.. Machine learning
1687 parameterization of the surface layer: bridging the observation-modeling gap. In *AGU Fall Meeting Abstracts*
1688 (Vol. 2019, pp. IN44A-04). 2019.

1689 Gagne, D. J., Chen, C. C., & Gettelman, A. Emulation of bin Microphysical Processes with machine learning. In *100th*
1690 *American Meteorological Society Annual Meeting*. AMS. 2020.

1691 Gagne, D. J., Christensen, H. M., Subramanian, A. C., & Monahan, A. H. Machine learning for stochastic
1692 parameterization: Generative adversarial networks in the Lorenz'96 model. *J ADV MODEL EARTH SY*,
1693 *12*(3), e2019MS001896. <https://doi.org/10.1029/2019MS001896>. 2020.

1694 Garg, S., Rasp, S., & Thuerey, N. WeatherBench Probability: A benchmark dataset for probabilistic medium-range
1695 weather forecasting along with deep learning baseline models. arXiv preprint arXiv:2205.00865.
1696 <https://doi.org/10.48550/arXiv.2205.00865>. 2 May 2022.

1697 George, T., Gupta, A., & Sarin, V. A recommendation system for preconditioned iterative solvers. *IEEE DATA*
1698 *MINING* (pp. 803-808). IEEE. <https://doi.org/10.1109/ICDM.2008.105>. 2008.

1699 Gettelman, A., Gagne, D. J., Chen, C. C., Christensen, M. W., Lebo, Z. J., Morrison, H., & Gantos, G. Machine
1700 learning the warm rain process. *J ADV MODEL EARTH SY*, *13*(2), e2020MS002268.
1701 <https://doi.org/10.1029/2020MS002268>. 2021.

1702 Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., ... & Bengio, Y. Generative
1703 adversarial networks. *COMMUN ACM*, *63*(11), 139-144. <https://doi.org/10.1145/3422622>. 2020.

1704 Goodfellow, I., Yoshua B., & Aaron C. Deep learning. *MIT press*. 2016.

1705 Grigo, C., & Koutsourelakis, P. S. (2019). A physics-aware, probabilistic machine learning framework for coarse-
1706 graining high-dimensional systems in the Small Data regime. *J COMPUT PHYS*, *397*, 108842.
1707 <https://doi.org/10.1016/j.jcp.2019.05.053>. 2019.

1708 Gurvan, M., Bourdallé-Badie, R., Chanut, J., Clementi, E., Coward, A., Ethé, C., ... & Samson, G. NEMO ocean
1709 engine, Institut Pierre-Simon Laplace (IPSL), Zenodo. 2019.

1710 Ham, Y. G., Kim, J. H., & Luo, J. J. Deep learning for multi-year ENSO forecasts. *Nature*, *573*(7775), 568-572.
1711 <https://doi.org/10.1038/s41586-019-1559-7>. 2019.

1712 Ham, Y. G., Kim, J. H., Kim, E. S., & On, K. W. Unified deep learning model for El Niño/Southern Oscillation
1713 forecasts by incorporating seasonality in climate data. *SCI BULL*, *66*(13), 1358-1366.
1714 <https://doi.org/10.1016/j.scib.2021.03.009>. 2021.

1715 Han, Y., Zhang, G. J., Huang, X., & Wang, Y. A moist physics parameterization based on deep learning. *J ADV*
1716 *MODEL EARTH SY*, *12*(9), e2020MS002076. <https://doi.org/10.1029/2020MS002076>. 2020.

1717 Harder, P., Watson-Parris, D., Stier, P., Strassel, D., Gauger, N. R., & Keuper, J. Physics-informed learning of aerosol
1718 microphysics. *Environmental Data Science*, *1*, e20. <https://doi.org/10.1017/eds.2022.22>. 2022.

1719 Harris, L., Chen, X., Putman, W., Zhou, L., & Chen, J. H. A scientific description of the GFDL finite-volume cubed-
1720 sphere dynamical core. <https://doi.org/10.25923/6nhs-5897>. 2021.

1721 Hastie, T., Tibshirani, R., Friedman, J. H., & Friedman, J. H. *The elements of statistical learning: data mining,*
1722 *inference, and prediction* (Vol. 2, pp. 1-758). New York: springer. 2009.

1723 He, K., Zhang, X., Ren, S., & Sun, J. Deep residual learning for image recognition. *PROC CVPR IEEE* (pp. 770-778).
1724 2016.

1725 He, X., Liu, S., Xu, T., Yu, K., Gentine, P., Zhang, Z., ... & Wu, D. Improving predictions of evapotranspiration by
1726 integrating multi-source observations and land surface model. *AGR WATER MANAGE*, 272, 107827.
1727 <https://doi.org/10.1016/j.agwat.2022.107827>. 2022.

1728 Hewamalage, H., Ackermann, K., & Bergmeir, C. Forecast Evaluation for Data Scientists: Common Pitfalls and Best
1729 Practices. *arXiv preprint arXiv:2203.10716*. <https://doi.org/10.48550/arXiv.2203.10716>. 21 March 2022.

1730 Hochreiter, S., & Schmidhuber, J. LSTM can solve hard long time lag problems. *ADV NEUR IN*, 9. 1996.

1731 Holloway, A., & Chen, T. Y. Neural networks for predicting the behavior of preconditioned iterative solvers. In
1732 *International Conference on Computational Science* (pp. 302-309). Springer, Berlin, Heidelberg.
1733 https://doi.org/10.1007/978-3-540-72584-8_39. 2007.

1734 Hornik, K., Stinchcombe, M., & White, H. Multilayer feedforward networks are universal approximators. *Neural
1735 networks*, 2(5), 359-366. [https://doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8). 1989.

1736 Horvat, C., & Roach, L. A. WIFF1. 0: a hybrid machine-learning-based parameterization of wave-induced sea ice floe
1737 fracture. *GEOSCI MODEL DEV*, 15(2), 803-814. <https://doi.org/10.5194/gmd-15-803-2022>. 2022

1738 Hsieh, W. W. Introduction to Environmental Data Science. Cambridge University Press. 2023.

1739 Hu, Y., Chen, L., Wang, Z., & Li, H. SwinVRNN: A Data-Driven Ensemble Forecasting Model via Learned
1740 Distribution Perturbation. *J ADV MODEL EARTH SY*, 15(2), e2022MS003211.
1741 <https://doi.org/10.1029/2022MS003211>. 2023.

1742 Huang, Z., England, M., Davenport, J. H., & Paulson, L. C. Using machine learning to decide when to precondition
1743 cylindrical algebraic decomposition with Groebner bases. In 2016 18th INT SYMP SYMB NUMERI
1744 (SYNASC) (pp. 45-52). IEEE. <https://doi.org/10.1109/SYNASC.2016.020>. 2016.

1745 Johnson, S. J., Stockdale, T. N., Ferranti, L., Balmaseda, M. A., Molteni, F., Magnusson, L., ... & Monge-Sanz, B. M.
1746 SEAS5: the new ECMWF seasonal forecast system. *GEOSCI MODEL DEV*, 12(3), 1087-1117.
1747 <https://doi.org/10.5194/gmd-12-1087-2019>. 2019.

1748 Kapp-Schworer, L., Graubner, A., Kim, S., & Kashinath, K. Spatio-temporal segmentation and tracking of weather
1749 patterns with light-weight Neural Networks. 2020.

1750 Karpatne, A., Atluri, G., Faghmous, J. H., Steinbach, M., Banerjee, A., Ganguly, A., ... & Kumar, V. Theory-guided
1751 data science: A new paradigm for scientific discovery from data. *IEEE T KNOWL DATA EN*, 29(10), 2318-
1752 2331. <https://doi.org/10.1109/TKDE.2017.2720168>. 2017.

1753 Karunasinghe, D. S., & Liang, S. Y. Chaotic time series prediction with a global model: Artificial neural network. *J
1754 HYDROL*, 323(1-4), 92-105. <https://doi.org/10.1016/j.jhydrol.2005.07.048>. 2006.

1755 Kashinath, K., Mustafa, M., Albert, A., Wu, J. L., Jiang, C., Esmaeilzadeh, S., ... & Prabhat. Physics-informed machine
1756 learning: case studies for weather and climate modelling. *Philosophical Transactions of the Royal Society A*,
1757 379(2194), 20200093. <https://doi.org/10.1098/rsta.2020.0093>. 2021.

1758 Keisler, R. Forecasting Global Weather with Graph Neural Networks. *arXiv preprint arXiv:2202.07575*.
1759 <https://doi.org/10.48550/arXiv.2202.07575>. 15 February 2022.

1760 Kelotra, A., & Pandey, P. (2020). Stock market prediction using optimized deep-convlstm model. *Big Data*, 8(1), 5-
1761 24. <https://doi.org/10.48550/arXiv.2202.07575>. 11 February 2022.

1762 Kim, J., Kwon, M., Kim, S. D., Kug, J. S., Ryu, J. G., & Kim, J. Spatiotemporal neural network with attention
1763 mechanism for El Niño forecasts. *SCI REP-UK*, 12(1), 1-15. <https://doi.org/10.1038/s41598-022-10839-z>.
1764 2022.

1765 Kipf, T. N., & Welling, M. Semi-supervised classification with graph convolutional networks. arXiv preprint
1766 arXiv:1609.02907. <https://doi.org/10.48550/arXiv.1609.02907>. 9 September 2016.

1767 Kochkov, D., Smith, J. A., Alieva, A., Wang, Q., Brenner, M. P., & Hoyer, S. Machine learning–accelerated
1768 computational fluid dynamics. *Proceedings of the National Academy of Sciences*, 118(21), e2101784118.
1769 <https://doi.org/10.1073/pnas.2101784118>. 2021.

1770 Krasnopolsky, V. M., Chalikov, D. V., & Tolman, H. L. A neural network technique to improve computational
1771 efficiency of numerical oceanic models. *OCEAN MODEL*, 4(3-4), 363-383. [https://doi.org/10.1016/S1463-
1772 5003\(02\)00010-0](https://doi.org/10.1016/S1463-5003(02)00010-0). 2002.

1773 Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Chalikov, D. V. New approach to calculation of atmospheric model
1774 physics: Accurate and fast neural network emulation of longwave radiation in a climate model. *MON
1775 WEATHER REV*, 133(5), 1370-1383. <https://doi.org/10.1175/MWR2923.1>. 2005.

1776 Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Belochitski, A. A. Using ensemble of neural networks to learn
1777 stochastic convection parameterizations for climate and numerical weather prediction models from data
1778 simulated by a cloud resolving model. *Advances in Artificial Neural Systems*, 2013.
1779 <https://doi.org/10.1155/2013/485913>. 2013.

1780 Kuefler, E., & Chen, T. Y. On using reinforcement learning to solve sparse linear systems. In *International Conference
1781 on Computational Science* (pp. 955-964). Springer, Berlin, Heidelberg. [https://doi.org/10.1007/978-3-540-
1782 69384-0_100](https://doi.org/10.1007/978-3-540-69384-0_100). 2008.

1783 Ladický, L. U., Jeong, S., Solenthaler, B., Pollefeys, M., & Gross, M. Data-driven fluid simulations using regression
1784 forests. *ACM T GRAPHIC (TOG)*, 34(6), 1-9. <https://doi.org/10.1145/2816795.2818129>. 2015.

1785 Lam, R., Sanchez-Gonzalez, A., Willson, M., Wirnsberger, P., Fortunato, M., Pritzel, A., ... & Battaglia, P. GraphCast:
1786 Learning skillful medium-range global weather forecasting. *arXiv preprint arXiv:2212.12794*.
1787 <https://doi.org/10.48550/arXiv.2212.12794>. 24 December 2022.

1788 Lanthaler, S., Mishra, S., & Karniadakis, G. E. Error estimates for deepnets: A deep learning framework in infinite
1789 dimensions. *Transactions of Mathematics and Its Applications*, 6(1), tnac001.
1790 <https://doi.org/10.1093/imatrm/tnac001>. 2022.

1791 LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. Gradient-based learning applied to document recognition. *P IEEE*,
1792 86(11), 2278-2324. <https://doi.org/10.1109/5.726791>. 1998.

1793 Leufen, L. H., & Schädler, G. Calculating the turbulent fluxes in the atmospheric surface layer with neural networks.
1794 *GEOSCI MODEL DEV*, 12(5), 2033-2047. <https://doi.org/10.5194/gmd-12-2033-2019>. 2019.

1795 Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Stuart, A., Bhattacharya, K., & Anandkumar, A. Multipole graph
1796 neural operator for parametric partial differential equations. *ADV NEUR IN*, 33, 6755-6766. 2020a.

1797 Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Bhattacharya, K., Stuart, A., & Anandkumar, A. Neural operator:
1798 Graph kernel network for partial differential equations. *arXiv preprint arXiv:2003.03485*.
1799 <https://doi.org/10.48550/arXiv.2003.03485>. 7 March 2020.

1800 Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Bhattacharya, K., Stuart, A., & Anandkumar, A. (2020). Fourier
1801 neural operator for parametric partial differential equations. *arXiv preprint arXiv:2010.08895*.
1802 <https://doi.org/10.48550/arXiv.2010.08895>. 18 October 2020.

1803 Lopez-Gomez, I., McGovern, A., Agrawal, S., & Hickey, J. Global extreme heat forecasting using neural weather
1804 models. *Artificial Intelligence for the Earth Systems*, 2(1), e220035. [https://doi.org/10.1175/AIES-D-22-](https://doi.org/10.1175/AIES-D-22-0035.1)
1805 [0035.1](https://doi.org/10.1175/AIES-D-22-0035.1). 2023.

1806 Lorenz, E. N. Predictability: A problem partly solved, in *Proceedings of Seminar on Predictability*, 4–8 September
1807 1995. <https://doi.org/10.1017/CBO9780511617652.004>. 1995

1808 Lu, L., Jin, P., & Karniadakis, G. E. DeepoNet: Learning nonlinear operators for identifying differential equations
1809 based on the universal approximation theorem of operators. *arXiv preprint arXiv:1910.03193*.
1810 <https://doi.org/10.48550/arXiv.1910.03193>. 8 October 2019.

1811 Lundberg S., & Lee S. A Unified Approach to Interpreting Model Predictions. *ADV NEUR IN*, 30, 4768–4777. 2017.

1812 McCulloch, W. S., & Pitts, W. A logical calculus of the ideas immanent in nervous activity. *The B MATH BIOPHYS*,
1813 5(4), 115-133. <https://doi.org/10.1007/BF02478259>. 1943.

1814 McGovern, A., Lagerquist, R., Gagne, D. J., Jergensen, G. E., Elmore, K. L., Homeyer, C. R., & Smith, T. Making
1815 the black box more transparent: Understanding the physical implications of machine learning. *B AM*
1816 *METEOROL SOC*, 100(11), 2175-2199. <https://doi.org/10.1175/BAMS-D-18-0195.1>. 2019.

1817 Meyer, D., Hogan, R. J., Dueben, P. D., & Mason, S. L. Machine learning emulation of 3D cloud radiative effects. *J*
1818 *ADV MODEL EARTH SY*, 14(3), e2021MS002550. <https://doi.org/10.1029/2021MS002550>. 2022.

1819 Moishin, M., Deo, R. C., Prasad, R., Raj, N., & Abdulla, S. Designing deep-based learning flood forecast model with
1820 ConvLSTM hybrid algorithm. *IEEE ACCESS*, 9, 50982-50993.
1821 <https://doi.org/10.1109/ACCESS.2021.3065939>. 2021.

1822 Molina, M. J., O'Brien, T. A., Anderson, G., Ashfaq, M., Bennett, K. E., Collins, W. D., ... & Ullrich, P. A. A Review
1823 of Recent and Emerging Machine Learning Applications for Climate Variability and Weather Phenomena.
1824 *Artificial Intelligence for the Earth Systems*, 1-46. <https://doi.org/10.1175/AIES-D-22-0086.1>. 2023.

1825 Mooers, G., Pritchard, M., Beucler, T., Ott, J., Yacalis, G., Baldi, P., & Gentine, P. (2021). Assessing the Potential of
1826 Deep Learning for Emulating Cloud Superparameterization in Climate Models With Real-Geography
1827 Boundary Conditions. *J ADV MODEL EARTH SY*, 13(5), e2020MS002385.
1828 <https://doi.org/10.1029/2020MS002385>. 2021.

1829 Mudigonda, M., Kim, S., Mahesh, A., Kahou, S., Kashinath, K., Williams, D., ... & Prabhat, M. Segmenting and
1830 tracking extreme climate events using neural networks. In *Deep Learning for Physical Sciences (DLPS)*
1831 *Workshop, held with NIPS Conference*. 2017.

1832 Nelsen, N. H., & Stuart, A. M. The random feature model for input-output maps between banach spaces. *SIAM J SCI*
1833 *COMPUT*, 43(5), A3212-A3243. <https://doi.org/10.1137/20M133957X>. 2021.

1834 Nguyen, T., Brandstetter, J., Kapoor, A., Gupta, J. K., & Grover, A. ClimaX: A foundation model for weather and
1835 climate. arXiv preprint arXiv:2301.10343. <https://doi.org/10.48550/arXiv.2301.10343>. 24 January 2023.

1836 Nielsen, A. H., Iosifidis, A., & Karstoft, H. Forecasting large-scale circulation regimes using deformable
1837 convolutional neural networks and global spatiotemporal climate data. SCI REP-UK, 12(1), 1-12.
1838 [https://doi.org/10.1175/1520-0469\(1995\)052%3C1237:WRRAS%3E2.0.CO;2](https://doi.org/10.1175/1520-0469(1995)052%3C1237:WRRAS%3E2.0.CO;2). 2022.

1839 O'Brien, T. A., Risser, M. D., Loring, B., Elbashandy, A. A., Krishnan, H., Johnson, J., ... & Collins, W. D. Detection
1840 of atmospheric rivers with inline uncertainty quantification: TECA-BARD v1. 0.1. GEOSCI MODEL DEV,
1841 13(12), 6131-6148. <https://doi.org/10.5194/gmd-13-6131-2020>. 2020.

1842 O'Gorman, P. A., & Dwyer, J. G. Using machine learning to parameterize moist convection: Potential for modeling
1843 of climate, climate change, and extreme events. J ADV MODEL EARTH SY, 10(10), 2548-2563.
1844 <https://doi.org/10.1029/2018MS001351>. 2018.

1845 O'Leary, J., Paulson, J. A., & Mesbah, A. Stochastic physics-informed neural ordinary differential equations. J
1846 COMPUT PHYS, 468, 111466. <https://doi.org/10.1016/j.jcp.2022.111466>. 2022.

1847 Ott, J., Pritchard, M., Best, N., Linstead, E., Curcic, M., & Baldi, P. A Fortran-Keras deep learning bridge for scientific
1848 computing. *Scientific Programming*, 2020. <https://doi.org/10.1155/2020/8888811>. 2020.

1849 Pal, S., & Sharma, P. A review of machine learning applications in land surface modeling. Earth, 2(1), 174-190.
1850 <https://doi.org/10.3390/earth2010011>. 2021.

1851 Palmer, T. A vision for numerical weather prediction in 2030. *arXiv preprint arXiv:2007.04830*.
1852 <https://doi.org/10.48550/arXiv.2007.04830>. 3 July 2020.

1853 Pan, S., Pan, N., Tian, H., Friedlingstein, P., Sitch, S., Shi, H., ... & Running, S. W. Evaluation of global terrestrial
1854 evapotranspiration using state-of-the-art approaches in remote sensing, machine learning and land surface
1855 modeling. HYDROL EARTH SYST SC, 24(3), 1485-1509. <https://doi.org/10.5194/hess-24-1485-2020>.
1856 2020.

1857 Patel, R. G., Trask, N. A., Wood, M. A., & Cyr, E. C. A physics-informed operator regression framework for extracting
1858 data-driven continuum models. COMPUT METHOD APPL M, 373, 113500.
1859 <https://doi.org/10.1016/j.cma.2020.113500>. 2021.

1860 Pathak, J., Subramanian, S., Harrington, P., Raja, S., Chattopadhyay, A., Mardani, M., ... & Anandkumar, A. (2022).
1861 Fourcastnet: A global data-driven high-resolution weather model using adaptive fourier neural operators.
1862 *arXiv preprint arXiv:2202.11214*. <https://doi.org/10.48550/arXiv.2202.11214>. 22 February 2022.

1863 Peairs, L., & Chen, T. Y. Using reinforcement learning to vary the m in GMRES (m). PROCEDIA COMPUT SCI, 4,
1864 2257-2266. <https://doi.org/10.1016/j.procs.2011.04.246>. 2011.

1865 Pelissier, C., Frame, J., & Nearing, G. Combining parametric land surface models with machine learning. INT
1866 GEOSCI REMOTE SE, (pp. 3668-3671). IEEE. <https://doi.org/10.1109/IGARSS39084.2020.9324607>.
1867 2020.

1868 Pincus, R., Mlawer, E. J., & Delamere, J. S. (2019). Balancing accuracy, efficiency, and flexibility in radiation
1869 calculations for dynamical models. J ADV MODEL EARTH SY, 11(10), 3074-3089.
1870 <https://doi.org/10.1029/2019MS001621>. 2019.

1871 Prabhat, P., Kashinath, K., Mudigonda, M., Kim, S., Kapp-Schwoerer, L., Graubner, A., & Collins, W. ClimateNet:
1872 An expert-labelled open dataset and Deep Learning architecture for enabling high-precision analyses of
1873 extreme weather. *GEOSCI MODEL DEV*, 14(1), 107-124. <https://doi.org/10.5194/gmd-14-107-2021>. 2021.

1874 Psaros, A. F., Meng, X., Zou, Z., Guo, L., & Karniadakis, G. E. Uncertainty quantification in scientific machine
1875 learning: Methods, metrics, and comparisons. *arXiv preprint arXiv:2201.07766*.
1876 <https://doi.org/10.48550/arXiv.2201.07766>. 19 January 2022.

1877 Rasp, S. Coupled online learning as a way to tackle instabilities and biases in neural network parameterizations:
1878 general algorithms and Lorenz 96 case study (v1. 0). *GEOSCI MODEL DEV*, 13(5), 2185-2196.
1879 <https://doi.org/10.48550/arXiv.1907.01351>. 24 March 2020.

1880 Rasp, S., Dueben, P. D., Scher, S., Weyn, J. A., Mouatadid, S., & Thuerey, N. WeatherBench: a benchmark data set
1881 for data-driven weather forecasting. *J ADV MODEL EARTH SY*, 12(11), e2020MS002203.
1882 <https://doi.org/10.1029/2020MS002203>. 2020.

1883 Rasp, S., Pritchard, M. S., & Gentine, P. Deep learning to represent subgrid processes in climate models. *Proceedings*
1884 *of the National Academy of Sciences*, 115(39), 9684-9689. <https://doi.org/10.1073/pnas.1810286115>. 2018.

1885 Rasp, S., & Thuerey, N. Data-driven medium-range weather prediction with a resnet pretrained on climate simulations:
1886 A new model for weatherbench. *J ADV MODEL EARTH SY*, 13(2), e2020MS002405.
1887 <https://doi.org/10.1029/2020MS002405>. 2021.

1888 Ravuri, S., Lenc, K., Willson, M., Kangin, D., Lam, R., Mirowski, P., ... & Mohamed, S. Skilful precipitation
1889 nowcasting using deep generative models of radar. *Nature*, 597(7878), 672-677.
1890 <https://doi.org/10.1038/s41586-021-03854-z>. 2021.

1891 Rizzuti, G., Siahkoobi, A., & Herrmann, F. J. Learned iterative solvers for the Helmholtz equation. In *81st EAGE*
1892 *Conference and Exhibition 2019* (Vol. 2019, No. 1, pp. 1-5). European Association of Geoscientists &
1893 Engineers. <https://doi.org/10.3997/2214-4609.201901542>. 2019.

1894 Ronneberger, O., Fischer, P., & Brox, T. U-net: Convolutional networks for biomedical image segmentation. In
1895 *International Conference on Medical image computing and computer-assisted intervention* (pp. 234-241).
1896 Springer, Cham. https://doi.org/10.1007/978-3-319-24574-4_28. 2015.

1897 Rosier, S. H., Bull, C., Woo, W. L., & Gudmundsson, G. H. Predicting ocean-induced ice-shelf melt rates using deep
1898 learning. *The Cryosphere*, 17(2), 499-518. <https://doi.org/10.5194/tc-17-499-2023>. 2023.

1899 Ross, A., Li, Z., Perezhogin, P., Fernandez-Granda, C., & Zanna, L. Benchmarking of machine learning ocean subgrid
1900 parameterizations in an idealized model. *J ADV MODEL EARTH SY*, 15(1), e2022MS003258.
1901 <https://doi.org/10.1029/2022MS003258>. 2023.

1902 Rupe, A., Kashinath, K., Kumar, N., & Crutchfield, J. P. (2023). Physics-Informed Representation Learning for
1903 Emergent Organization in Complex Dynamical Systems. *arXiv preprint arXiv:2304.12586*.
1904 <https://doi.org/10.48550/arXiv.2304.12586>. 25 April 2023.

1905 Russell S. & Norvig P. *Artificial Intelligence: A Modern Approach* (Fourth Global Edition). Pearson Education. 2021.

1906 Samek, W., Montavon, G., Lapuschkin, S., Anders, C. J., & Müller, K. R. Explaining deep neural networks and
1907 beyond: A review of methods and applications. *P IEEE*, 109(3), 247-278.
1908 <https://doi.org/10.1109/JPROC.2021.3060483>. 2021.

1909 Sawada, Y. Machine learning accelerates parameter optimization and uncertainty assessment of a land surface model.
1910 *J GEOPHYS RES-ATMOS*, 125(20), e2020JD032688. <https://doi.org/10.1029/2020JD032688>. 2020.

1911 Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M., & Monfardini, G. The graph neural network model. *IEEE T*
1912 *NEURAL NETWOR*, 20(1), 61-80. <https://doi.org/10.1109/TNN.2008.2005605>. 2008.

1913 Scher, S. Toward data-driven weather and climate forecasting: Approximating a simple general circulation model with
1914 deep learning. *GEOPHYS RES LETT*, 45(22), 12-616. <https://doi.org/10.1029/2018GL080704>. 2018.

1915 Scher, S., & Messori, G. Weather and climate forecasting with neural networks: using general circulation models
1916 (GCMs) with different complexity as a study ground. *GEOSCI MODEL DEV*, 12(7), 2797-2809.
1917 <https://doi.org/10.5194/gmd-12-2797-2019>. 2019.

1918 Shi, X., Chen, Z., Wang, H., Yeung, D. Y., Wong, W. K., & Woo, W. C. Convolutional LSTM network: A machine
1919 learning approach for precipitation nowcasting. *ADV NEUR IN*, 28. 2015.

1920 Taylor, J., & Feng, M. A Deep Learning Model for Forecasting Global Monthly Mean Sea Surface Temperature
1921 Anomalies. *arXiv preprint arXiv:2202.09967*. <https://doi.org/10.48550/arXiv.2202.09967>. 21 February
1922 2022.

1923 Tibshirani, R., & Friedman, J. H. *The elements of statistical learning [electronic resource]: data mining, inference,*
1924 *and prediction: with 200 full-color illustrations*. Springer. 2001.

1925 Tompson, J., Schlachter, K., Sprechmann, P., & Perlin, K. Accelerating eulerian fluid simulation with convolutional
1926 networks. In *International Conference on Machine Learning* (pp. 3424-3433). PMLR. 2017.

1927 Toms, B. A., Barnes, E. A., & Ebert-Uphoff, I. Physically interpretable neural networks for the geosciences:
1928 Applications to earth system variability. *J ADV MODEL EARTH SY*, 12(9), e2019MS002002.
1929 <https://doi.org/10.1029/2019MS002002>. 2020.

1930 Turing, A. M., Computing Machinery and Intelligence, *Mind*, Volume LIX, Issue 236, Pages 433–460,
1931 <https://doi.org/10.1093/mind/LIX.236.433>. 1950.

1932 Ukkonen, P., & Mäkelä, A. Evaluation of machine learning classifiers for predicting deep convection. *J ADV MODEL*
1933 *EARTH SY*, 11(6), 1784-1802. <https://doi.org/10.1029/2018MS001561>. 2019.

1934 Ukkonen, P., Pincus, R., Hogan, R. J., Pagh Nielsen, K., & Kaas, E. (2020). Accelerating radiation computations for
1935 dynamical models with targeted machine learning and code optimization. *J ADV MODEL EARTH SY*,
1936 12(12), e2020MS002226. <https://doi.org/10.1029/2020MS002226>. 2020.

1937 United Nations Educational, Scientific and Cultural Organization. Recommendations on the Ethics of Artificial
1938 Intelligence. 2021.

1939 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. Attention is all you
1940 need. *ADV NEUR IN*, 30. 2017.

1941 Virnodkar, S. S., Pachghare, V. K., Patil, V. C., & Jha, S. K. Remote sensing and machine learning for crop water
1942 stress determination in various crops: a critical review. *PRECIS AGRIC*, 21(5), 1121-1155.
1943 <https://doi.org/10.1007/s11119-020-09711-9>. 2020.

1944 Vlachas, P. R., Byeon, W., Wan, Z. Y., Sapsis, T. P., & Koumoutsakos, P. Data-driven forecasting of high-dimensional
1945 chaotic systems with long short-term memory networks. *Proceedings of the Royal Society A: Mathematical,*
1946 *Physical and Engineering Sciences*, 474(2213), 20170844. <https://doi.org/10.1098/rspa.2017.0844>. 2018.

1947 Walters, D., Boutle, I., Brooks, M., Melvin, T., Stratton, R., Vosper, S., ... & Xavier, P. The Met Office unified model
1948 global atmosphere 6.0/6.1 and JULES global land 6.0/6.1 configurations. *GEOSCI MODEL DEV*, 10(4),
1949 1487-1520. <https://doi.org/10.5194/gmd-10-1487-2017>. 2017.

1950 Wang, S., Wang, H., & Perdikaris, P. Learning the solution operator of parametric partial differential equations with
1951 physics-informed DeepONets. *Science advances*, 7(40), eabi8605. <https://doi.org/10.1126/sciadv.abi8605>.
1952 2021.

1953 Wang, X., Han, Y., Xue, W., Yang, G., & Zhang, G. J. Stable climate simulations using a realistic general circulation
1954 model with neural network parameterizations for atmospheric moist physics and radiation processes.
1955 *GEOSCI MODEL DEV*, 15(9), 3923-3940. <https://doi.org/10.5194/gmd-15-3923-2022>. 2022.

1956 Wang, S., Sankaran, S., & Perdikaris, P. (2022b). Respecting causality is all you need for training physics-informed
1957 neural networks. *arXiv preprint arXiv:2203.07404*. <https://doi.org/10.48550/arXiv.2203.07404>. 14 March
1958 2022.

1959 Watson, P. A. Machine learning applications for weather and climate need greater focus on extremes. *ENVIRON RES*
1960 *LETT*, 17(11), 111004. <https://doi.org/10.1088/1748-9326/ac9d4e>. 2022.

1961 Watt-Meyer, O., Brenowitz, N. D., Clark, S. K., Henn, B., Kwa, A., McGibbon, J., ... & Bretherton, C. S. Correcting
1962 weather and climate models by machine learning nudged historical simulations. *GEOPHYS RES LETT*,
1963 48(15), e2021GL092555. <https://doi.org/10.1029/2021GL092555>. 2021.

1964 Watson-Parris, D. Machine learning for weather and climate are worlds apart. *Philosophical Transactions of the Royal*
1965 *Society A*, 379(2194), 20200098. <https://doi.org/10.1098/rsta.2020.0098>. 2021.

1966 Watson-Parris, D., Rao, Y., Olivić, D., Seland, Ø., Nowack, P., Camps-Valls, G., ... & Roesch, C. ClimateBench v1.
1967 0: A Benchmark for Data-Driven Climate Projections. *J ADV MODEL EARTH SY*, 14(10),
1968 e2021MS002954. <https://doi.org/10.1029/2021MS002954>. 2022.

1969 Werbos, P. Beyond regression: new tools for prediction and analysis in the behavioral sciences. *Ph. D. dissertation,*
1970 *Harvard University*. 1974.

1971 Werbos, P. J. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10), 1550-
1972 1560. <https://doi.org/10.1109/5.58337>. 1990.

1973 Weyn, J. A., Durran, D. R., & Caruana, R. (2019). Can machines learn to predict weather? Using deep learning to
1974 predict gridded 500-hPa geopotential height from historical weather data. *J ADV MODEL EARTH SY*,
1975 11(8), 2680-2693. <https://doi.org/10.1029/2019MS001705>. 2019.

1976 Weyn, J. A., Durran, D. R., & Caruana, R. Improving data-driven global weather prediction using deep convolutional
1977 neural networks on a cubed sphere. *J ADV MODEL EARTH SY*, 12(9), e2020MS002109.
1978 <https://doi.org/10.1029/2020MS002109>. 2020.

1979 Weyn, J. A., Durran, D. R., Caruana, R., & Cresswell-Clay, N. Sub-seasonal forecasting with a large ensemble of
1980 deep-learning weather prediction models. *J ADV MODEL EARTH SY*, 13(7).
1981 <https://doi.org/10.1029/2021MS002502>. 2021.

1982 Wikner, A., Pathak, J., Hunt, B., Girvan, M., Arcomano, T., Szunyogh, I., ... & Ott, E. Combining machine learning
1983 with knowledge-based modeling for scalable forecasting and subgrid-scale closure of large, complex,
1984 spatiotemporal systems. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 30(5), 053111.
1985 <https://doi.org/10.1063/5.0005541>. 2020.

1986 Wu, K., & Xiu, D. Data-driven deep learning of partial differential equations in modal space. *J COMPUT PHYS*, 408,
1987 109307. <https://doi.org/10.1016/j.jcp.2020.109307>. 2020.

1988 Yamada, K., Katagiri, T., Takizawa, H., Minami, K., Yokokawa, M., Nagai, T., & Ogino, M. Preconditioner auto-
1989 tuning using deep learning for sparse iterative algorithms. In *2018 Sixth International Symposium on*
1990 *Computing and Networking Workshops (CANDARW)* (pp. 257-262). IEEE.
1991 <https://doi.org/10.1016/j.jcp.2020.109307>. 2018.

1992 Yang, C., Yang, X., & Xiao, X. Data-driven projection method in fluid simulation. *COMPUT ANIMAT VIRT W*,
1993 27(3-4), 415-424. <https://doi.org/10.1002/cav.1695>. 2016.

1994 Yeo, K., Grullon, D. E., Sun, F. K., Boning, D. S., & Kalagnanam, J. R. Variational inference formulation for a model-
1995 free simulation of a dynamical system with unknown parameters by a recurrent neural network. *SIAM J SCI*
1996 *COMPUT*, 43(2), A1305-A1335. <https://doi.org/10.1137/20M1323151>. 2021.

1997 Yuan, Z., Zhou, X., & Yang, T. Hetero-convlstm: A deep learning approach to traffic accident prediction on
1998 heterogeneous spatio-temporal data. In *Proceedings of the 24th ACM SIGKDD International Conference on*
1999 *Knowledge Discovery & Data Mining* (pp. 984-992). <https://doi.org/10.1145/3219819.3219922>. 2018.

2000 Yuval, J., & O’Gorman, P. A. (2020). Stable machine-learning parameterization of subgrid processes for climate
2001 modeling at a range of resolutions. *NAT COMMUN*, 11(1), 1-10. [https://doi.org/10.1038/s41467-020-](https://doi.org/10.1038/s41467-020-17142-3)
2002 [17142-3](https://doi.org/10.1038/s41467-020-17142-3). 2020.

2003 Yuval, J., O’Gorman, P. A., & Hill, C. N. Use of neural networks for stable, accurate and physically consistent
2004 parameterization of subgrid atmospheric processes with good performance at reduced precision. *GEOPHYS*
2005 *RES LETT*, 48(6), e2020GL091363. <https://doi.org/10.1029/2020GL091363>. 2021.

2006 Zagoruyko, S., & Komodakis, N. Wide residual networks. arXiv preprint arXiv:1605.07146.
2007 <https://doi.org/10.48550/arXiv.1605.07146>. 23 May 2016.

2008 Zanna, L., & Bolton, T. Data-driven equation discovery of ocean mesoscale closures. *GEOPHYS RES LETT*, 47(17),
2009 e2020GL088376. <https://doi.org/10.1029/2020GL088376>. 2020.

2010 Zhang, N., Zhou, X., Kang, M., Hu, B. G., Heuvelink, E., & Marcelis, L. F. Machine learning versus crop growth
2011 models: an ally, not a rival. *AOB PLANTS*, 15(2), plac061. <https://doi.org/10.1093/aobpla/plac061>. 2023.

2012 Zhao, W. L., Gentine, P., Reichstein, M., Zhang, Y., Zhou, S., Wen, Y., ... & Qiu, G. Y. Physics-constrained machine
2013 learning of evapotranspiration. *GEOPHYS RES LETT*, 46(24), 14496-14507.
2014 <https://doi.org/10.1029/2019GL085291>. 2019.

2015 Zhong, X., Ma, Z., Yao, Y., Xu, L., Wu, Y., & Wang, Z. WRF-ML v1. 0: a bridge between WRF v4. 3 and machine
2016 learning parameterizations and its application to atmospheric radiative transfer. *GEOSCI MODEL DEV*,
2017 16(1), 199-209. <https://doi.org/10.5194/gmd-16-199-2023>. 2023.

2018 Zhou, L., Lin, S. J., Chen, J. H., Harris, L. M., Chen, X., & Rees, S. L. Toward convective-scale prediction within the
2019 next generation global prediction system. *B AM METEOROL SOC*, 100(7), 1225-1243.
2020 <https://doi.org/10.1175/BAMS-D-17-0246.1>. 2019.

2021