

Data-Driven Equation Discovery of a Cloud Cover Parameterization

Arthur Grundner^{1,2}, Tom Beucler³, Pierre Gentine², and Veronika Eyring^{1,4}

¹Deutsches Zentrum für Luft- und Raumfahrt e.V. (DLR), Institut für Physik der Atmosphäre,
Oberpfaffenhofen, Germany

²Center for Learning the Earth with Artificial Intelligence And Physics (LEAP), Columbia University,
New York, NY, USA

³Institute of Earth Surface Dynamics, University of Lausanne, Lausanne, Switzerland

⁴University of Bremen, Institute of Environmental Physics (IUP), Bremen, Germany

Key Points:

- We systematically derive and evaluate cloud cover parameterizations of various complexity from global storm-resolving simulation output
- Using symbolic regression combined with physical constraints, we find a new interpretable equation balancing performance and simplicity
- Our data-driven cloud cover equation can be retuned with few samples, facilitating transfer learning to generalize to other realistic data

Abstract

A promising method for improving the representation of clouds in climate models, and hence climate projections, is to develop machine learning-based parameterizations using output from global storm-resolving models. While neural networks can achieve state-of-the-art performance, they are typically climate model-specific, require post-hoc tools for interpretation, and struggle to predict outside of their training distribution. To avoid these limitations, we combine symbolic regression, sequential feature selection, and physical constraints in a hierarchical modeling framework. This framework allows us to discover new equations diagnosing cloud cover from coarse-grained variables of global storm-resolving model simulations. These analytical equations are interpretable by construction and easily transferable to other grids or climate models. Our best equation balances performance and complexity, achieving a performance comparable to that of neural networks ($R^2 = 0.94$) while remaining simple (with only 13 trainable parameters). It reproduces cloud cover distributions more accurately than the Xu-Randall scheme across all cloud regimes (Hellinger distances < 0.09), and matches neural networks in condensate-rich regimes. When applied and fine-tuned to the ERA5 reanalysis, the equation exhibits superior transferability to new data compared to all other optimal cloud cover schemes. Our findings demonstrate the effectiveness of symbolic regression in discovering interpretable, physically-consistent, and nonlinear equations to parameterize cloud cover.

Plain Language Summary

In climate models, cloud cover is usually expressed as a function of coarse, pixelated variables. Traditionally, this functional relationship is derived from physical assumptions. In contrast, machine learning approaches, such as neural networks, sacrifice interpretability for performance. In our approach, we use high-resolution climate model output to learn a hierarchy of cloud cover schemes from data. To bridge the gap between simple statistical methods and machine learning algorithms, we employ a symbolic regression method. Unlike classical regression, which requires providing a set of basis functions from which the equation is composed of, symbolic regression only requires mathematical operators (such as $+$, \times) that it learns to combine. By using a genetic algorithm, inspired by the process of natural selection, we discover an interpretable, nonlinear equation for cloud cover. This equation is simple, performs well, satisfies physical principles, and outperforms other algorithms when applied to new observationally-informed data.

1 Introduction

Due to computational constraints, climate models used to make future projections spanning multiple decades typically have horizontal resolutions of 50 – 100 km (Eyring et al., 2021). The coarse resolution necessitates the parameterization of many subgrid-scale processes (e.g., radiation, microphysics), which have a significant effect on model forecasts (Stensrud, 2009). Climate models, such as the state-of-the-art ICOSahedral Non-hydrostatic (ICON) model, exhibit long-standing systematic biases, especially related to cloud parameterizations (Crueger et al., 2018; Giorgetta et al., 2018). A fundamental component of the cloud parameterization package in ICON is its cloud cover scheme, which, in its current form, diagnoses fractional cloud cover from large-scale variables in every grid cell (Giorgetta et al., 2018; Mauritsen et al., 2019). As cloud cover is directly used in the radiation (Pincus & Stevens, 2013) and microphysics (Lohmann & Roeckner, 1996) parameterizations of ICON, its estimate directly influences the energy balance and the statistics of water vapor, cloud ice, and cloud water. The current cloud cover scheme in ICON, based on Sundqvist et al. (1989), nevertheless makes some crude empirical assumptions, such as a near-exclusive emphasis on relative humidity (see Grundner et al. (2022) for further discussion).

66 With the extended availability of high-fidelity data and increasingly sophisticated
 67 machine learning (ML) methods, ML algorithms have been developed for the parameter-
 68 ization of clouds and convection (e.g., Brenowitz and Bretherton (2018); Gentine et
 69 al. (2018); Krasnopolsky et al. (2013); O’Gorman and Dwyer (2018); see reviews by Beucler
 70 et al. (2022) and Gentine et al. (2021)). High-resolution atmospheric simulations on storm-
 71 resolving scales (horizontal resolutions of a few kilometers) resolve deep convective pro-
 72 cesses explicitly (Weisman et al., 1997), and provide useful training data with an improved
 73 physical representation of clouds and convection (Hohenegger et al., 2020; Stevens et al.,
 74 2020). There are only few approaches that learn parameterizations directly from obser-
 75 vations (e.g., McCandless et al. (2022)), as these are challenged by the sparsity and noise
 76 of observations (Rasp et al., 2018; Trenberth et al., 2009). Therefore, a two-step process
 77 might be required, in which the statistical model structure is first learned on high-resolution
 78 modeled data before its parameters are fine-tuned on observations (transfer learning),
 79 leveraging the advantage of the consistency of the modeled data for the initial training
 80 stage before having to deal with noisier observational data.

81 Neural networks and random forests have been routinely used for ML-based pa-
 82 rameterizations. Unlike traditional regression approaches, they are not limited to a par-
 83 ticular functional form provided by combining a set of basis functions. They are usually
 84 fast at inference time and can be trained with very little domain knowledge. However,
 85 this versatility comes at the cost of interpretability as explainable artificial intelligence
 86 (XAI) methods still face major challenges (Kumar et al., 2020; Molnar et al., 2021). Given
 87 this limitation, we ask: Can we create data-driven cloud cover schemes that are inter-
 88 pretable by construction without renouncing the high data fidelity of neural networks?

89 Here, we use a hierarchical modeling approach to systematically derive and evalu-
 90 ate a family of cloud cover (interpreted as the cloud area fraction) schemes, ranging from
 91 traditional physical (but semi-empirical) schemes and simple regression models to neu-
 92 ral networks. We evaluate them according to their Pareto optimality (i.e., whether they
 93 are the best performing model for their complexity). To bridge the gap between simple
 94 equations and high-performance neural networks, we apply equation discovery in a data-
 95 driven manner using state-of-the-art symbolic regression methods. In symbolic regres-
 96 sion, as opposed to regular regression, the user first specifies a set of mathematical op-
 97 erators instead of a set of basis functions. Based on these operators, the symbolic regres-
 98 sion library creates a random initial population of equations (Schmidt & Lipson, 2009).
 99 Inspired by the process of natural selection in the theory of evolution, symbolic regres-
 100 sion is usually implemented as a genetic algorithm that iteratively applies genetically mo-
 101 tivated operations (selection, crossover, mutation) to the set of candidate equations. At
 102 each step, the equations are ranked based on their performance and simplicity, so that
 103 the top equations can be selected to be included in the next population (Smits & Kotanchek,
 104 2005). Advantages of training/discovering analytical models instead of neural networks
 105 include an immediate view of model content (e.g., whether physical constraints are sat-
 106 isfied) and the ability to analyze the model structure directly using powerful mathemat-
 107 ical tools (e.g., perturbation theory, numerical stability analysis). Additionally, analyt-
 108 ical models are straightforward to communicate to the broader scientific community, to
 109 implement numerically, and fast to execute given the existence of optimized implemen-
 110 tations of well-known functions.

111 To our knowledge, Zanna and Bolton (2020) marks the first usage of automated,
 112 data-driven equation discovery for climate applications. Training on highly idealized data,
 113 they used a sparse regression technique called relevance vector machine to find an an-
 114 alytical model that parameterizes ocean eddies. In sparse regression, the user defines a
 115 library of terms, and the algorithm determines a linear combination of those terms that
 116 best matches the data while including as few terms as possible (Brunton et al., 2016; Rudy
 117 et al., 2017; Zhang & Lin, 2018; Champion et al., 2019). In a follow-up paper, Ross et
 118 al. (2023) employed symbolic regression to discover an improved equation, again trained

119 on idealized data, that performs similarly well as neural networks across various met-
 120 rics and has greater generalization capability. Nonetheless, they had to assume that the
 121 equation was linear in terms of its free/trainable parameters and additively separable
 122 as their method included an iterative approach to select suitable terms. For the selec-
 123 tion of terms, they took a human-in-the-loop approach rather than solely relying on the
 124 genetic algorithm. Additionally, the final discovered equation relied on high-order spa-
 125 tial derivatives, which may not be feasible to compute in a climate model. To prevent
 126 this issue, we only permit features we can either access or easily derive in the climate
 127 model.

128 Guiding questions for this study include: Using symbolic regression, can we auto-
 129 matically discover a physically consistent equation for cloud cover whose performance
 130 is competitive with that of neural networks? Given that modern symbolic regression li-
 131 braries can handle higher computational overhead, we want to relax prior assumptions
 132 of linearity or separability of the equation. Then, what can we learn about the cloud cover
 133 parameterization problem by sequentially selecting performance-maximizing features in
 134 different predictive models? Finally, how much better do simple models generalize and/or
 135 transfer to more realistic data sets?

136 We first introduce the data sets used for training, validation and testing (Sec 2),
 137 the diverse data-driven models used in this study (Sec 3), and evaluation metrics (Sec 4),
 138 before studying the feature rankings, performances and complexities of the different mod-
 139 els (Sec 5.1). We investigate their ability to reproduce cloud cover distributions (Sec 5.2)
 140 and adapt to the ERA5 reanalysis (Sec 5.3). We conclude with an analysis of the best
 141 analytical model we found using symbolic regression (Sec 6).

142 2 Data

143 In this section, we introduce the two data sets used to train and benchmark our
 144 cloud cover schemes: We first use storm-resolving ICON simulations to train high-fidelity
 145 models (Sec 2.1), before testing these models’ transferability to the ERA5 meteorolog-
 146 ical reanalysis, which is more directly informed by observations (Sec 2.2).

147 2.1 Global Storm-Resolving Model Simulations (DYAMOND)

148 As the source for our training data, we use output from global storm-resolving ICON
 149 simulations performed as part of the DYNAMICS of the Atmospheric general circulation
 150 Modeled On Non-hydrostatic Domains (DYAMOND) project. The project’s first phase
 151 (‘DYAMOND Summer’) included a simulation starting from August 1, 2018 (Stevens
 152 et al., 2019), while the second phase (‘DYAMOND Winter’) was initialized on January
 153 20, 2020 (Duras et al., 2021). In both phases, the ICON model simulated 40 days, pro-
 154 viding three-hourly output on a grid with a horizontal resolution of 2.47 km.

155 Following the methodology of Grundner et al. (2022), we coarse-grain the DYA-
 156 MOND data to an ICON grid with a typical climate model horizontal grid resolution of
 157 ≈ 80 km. Vertically, we coarse-grain the data from 58 to 27 layers below an altitude of
 158 21 km, which is the maximum altitude with clouds in the data set. For cloud cover, we
 159 first estimate the vertically maximal cloud cover values in each low-resolution grid cell
 160 before horizontally coarse-graining the resulting field. For all other variables, we take a
 161 three-dimensional integral over the high-resolution grid cells overlapping a given low-resolution
 162 grid cell. For details, we refer the reader to Appendix A of Grundner et al. (2022). Due
 163 to the sequential processing of some parameterization schemes in the ICON model, condensate-
 164 free clouds can occur in the simulation output. To instead ensure consistency between
 165 cloud cover and the other model variables, we follow Giorgetta et al. (2022) and man-
 166 ually set the cloud cover in the high-resolution grid cells to 100% when the cloud con-
 167 densate mixing ratio exceeds 10^{-6} kg/kg and to 0% otherwise.

168 We remove the first ten days of ‘DYAMOND Summer’ and ‘DYAMOND Winter’
 169 as spin-up, and discard columns that contain NaNs (3.15% of all columns). From the re-
 170 mainder, we keep a random subset of 28.5% of the data, while removing predominantly
 171 cloud-free cells to mitigate a class imbalance in the output (‘undersampling’ step). We
 172 then split the data into a training and a validation set, the latter of which is used for early
 173 stopping. To avoid high correlations between the training and validation sets, we divide
 174 the data set into six temporally connected parts. We choose the union of the second (\approx
 175 Aug 21 - Sept 1, 2016) and the fifth (\approx Feb 9 - Feb 19, 2020) part to create our valida-
 176 tion set. For all models except the traditional schemes, we additionally normalize mod-
 177 els’ features (or ‘inputs’) so that they have zero mean and unit variance on the train-
 178 ing set.

We define a set of 24 features \mathcal{F} that the models (discussed in Sec 3) can choose from. For clarity, we decompose \mathcal{F} into three subsets: $\mathcal{F} \stackrel{\text{def}}{=} \mathcal{F}_1 \cup \mathcal{F}_2 \cup \mathcal{F}_3$. The first subset, $\mathcal{F}_1 \stackrel{\text{def}}{=} \{U, q_v, q_c, q_i, T, p, \text{RH}\}$ groups the horizontal wind magnitude U [m/s] and thermodynamic variables known to influence cloud cover, namely specific humidity q_v [kg/kg], cloud water and ice mixing ratios q_c [kg/kg] and q_i [kg/kg], temperature T [K], pressure p [Pa], and relative humidity RH with respect to water, approximated as:

$$\text{RH} \approx 0.00263 \frac{p}{1\text{Pa}} q_v \exp \left[\frac{17.67(273.15\text{K} - T)}{T - 29.65\text{K}} \right]. \quad (1)$$

179 The second subset \mathcal{F}_2 contains the first and second vertical derivatives of all features in
 180 \mathcal{F}_1 . These derivatives are computed by fitting splines to every vertical profile of a given
 181 variable and differentiating the spline at the grid level heights to obtain derivatives on
 182 the irregular vertical grid. Finally, the third subset $\mathcal{F}_3 \stackrel{\text{def}}{=} \{z, \text{land}, p_s\}$ includes geo-
 183 metric height z [m] and the only two-dimensional variables, i.e., land fraction and sur-
 184 face pressure p_s [Pa].

185 In Grundner et al. (2022) we found it sufficient to diagnose cloud cover using in-
 186 formation from the close vertical neighborhood of a grid cell. By utilizing vertical deriva-
 187 tives to incorporate this information, we ensure the applicability of our cloud cover schemes
 188 to any vertical grid. Since our feature set \mathcal{F} contains all features appearing in our three
 189 baseline ‘traditional’ parameterizations (see Sec 3.1), we deem it comprehensive enough
 190 for the scope of our study.

191 2.2 Meteorological Reanalysis (ERA5)

192 To test the transferability of our cloud cover schemes to observational data, we also
 193 use the ERA5 meteorological reanalysis (Hersbach et al., 2018). We sample the first day
 194 of each quarter in 1979-2021 at a three-hourly resolution. The days from 2000-2006 are
 195 taken from ERA5.1, which uses an improved representation of the global-mean temper-
 196 atures in the upper troposphere and stratosphere. Depending on the ERA5 variable, they
 197 are either stored on an N320 reduced Gaussian (e.g., for cloud cover) or a T639 spec-
 198 tral (e.g., for temperature) grid. Using the CDO package (Schulzweida, 2019), we first
 199 remap all relevant variables to a regular Gaussian grid, and then to the unstructured ICON
 200 grid described in Sec 2.1. Vertically, we coarse-grain from approximately 90 to 27 lay-
 201 ers.

202 The univariate distributions of important features such as cloud water and ice do
 203 not match between the (coarse-grained) DYAMOND and (processed) ERA5 data. The
 204 maximal cloud ice values that are attained in the ERA5 data set are twice as large as
 205 in the DYAMOND data. We illustrate this in Fig 1, next to a comparison of the distri-
 206 butions of cloud water, relative humidity and temperature. Due to differences in the dis-
 207 tributions of cloud ice, cloud water and relative humidity, we consider our processed ERA5
 208 data a challenging data set to generalize to.

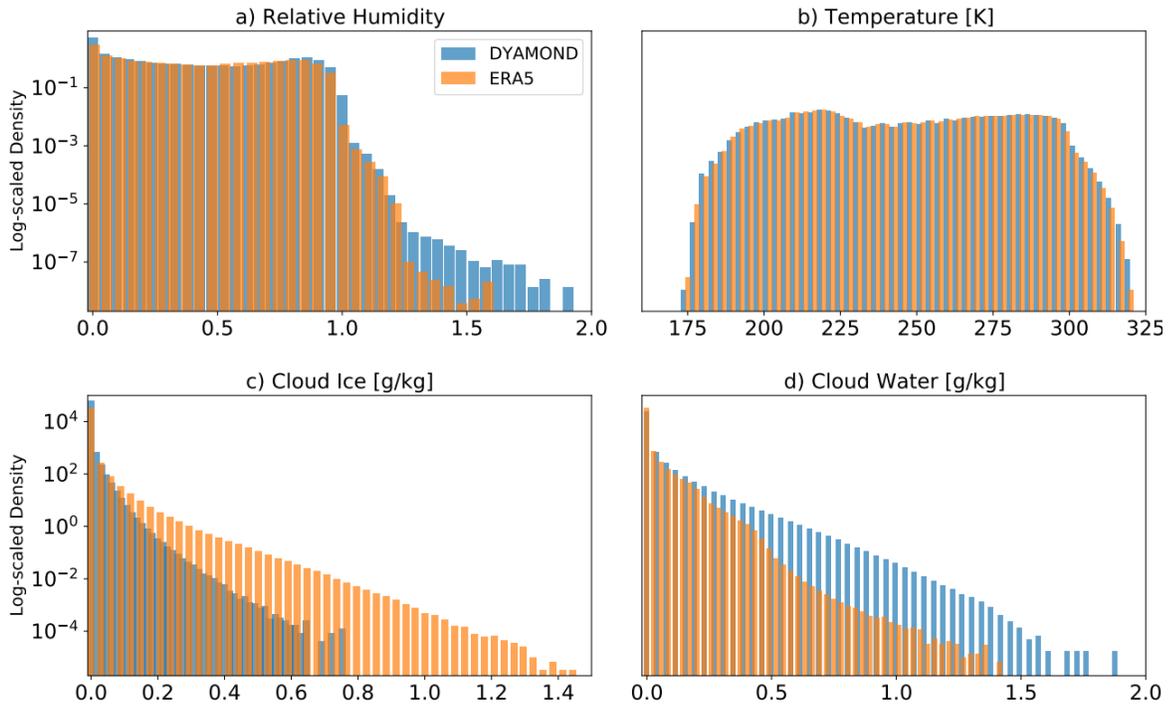


Figure 1. A comparison of the univariate distributions of four variables from the coarse-grained DYAMOND and ERA5 data sets. The y-axes are scaled logarithmically to visualize the distributions' tails. While cloud ice is often larger in our processed ERA5 data set, cloud water tends to be smaller than in the DYAMOND data. The distributions of temperature and relative humidity are comparable.

3 Data-Driven Modeling

We now introduce a family of data-driven cloud cover schemes. We adopt a hierarchical modeling approach and start with models that are interpretable by construction, i.e., linear models, polynomials, and traditional schemes. As a second step, we mostly focus on performance and therefore train deep neural networks (NNs) on the DYAMOND data. To bridge the gap between the best-performing and most interpretable models, we use symbolic regression to discover analytical cloud cover schemes from data. These schemes are complex enough to include relevant nonlinearities while remaining interpretable.

3.1 Existing Schemes

We first introduce three traditional diagnostic schemes for cloud cover and train them using the BFGS (Nocedal & Wright, 1999) and Nelder-Mead (Gao & Han, 2012) unconstrained optimizers (which outperform grid search methods in our case), each time choosing the model that minimizes the mean squared error (MSE) on the validation set. Before doing so, we multiply the output of each of the three schemes by 100 to obtain percent cloud cover values. The first is the Sundqvist scheme (Sundqvist et al., 1989), which is currently implemented in the ICON climate model (Giorgetta et al., 2018). The Sundqvist scheme expresses cloud cover as a monotonically increasing function of relative humidity. It assumes that cloud cover can only exist if relative humidity exceeds a critical relative humidity threshold RH_0 , which itself is a function of the fraction between surface pressure and pressure: If

$$\text{RH} > \text{RH}_0 \stackrel{\text{def}}{=} \text{RH}_{0,\text{top}} + (\text{RH}_{0,\text{surf}} - \text{RH}_{0,\text{top}}) \exp(1 - (p_s/p)^n), \quad (2)$$

then the Sundqvist cloud cover is given by

$$C_{\text{Sundqvist}} \stackrel{\text{def}}{=} 1 - \sqrt{\frac{\min\{\text{RH}, \text{RH}_{\text{sat}}\} - \text{RH}_{\text{sat}}}{\text{RH}_0 - \text{RH}_{\text{sat}}}}. \quad (3)$$

The Sundqvist scheme has four tunable parameters $\{\text{RH}_{0,\text{surf}}, \text{RH}_{0,\text{top}}, \text{RH}_{\text{sat}}, n\}$. As properly representing marine stratocumulus clouds in the Sundqvist scheme might require a different treatment (see e.g., Mauritsen et al. (2019)), we allow these parameters to differ between land and sea, which we separate using a land fraction threshold of 0.5.

The second scheme is a simplified version of the Xu-Randall scheme (Xu & Randall, 1996), which was found to outperform the Sundqvist scheme on CloudSat data (Wang et al., 2023). It additionally depends on cloud water and ice, ensuring that cloud cover is 0 in condensate-free grid cells. It can be formulated as

$$C_{\text{Xu-Randall}} \stackrel{\text{def}}{=} \min\{\text{RH}^\beta (1 - \exp(-\alpha(q_c + q_i))), 1\}. \quad (4)$$

The Xu-Randall scheme has only two tuning parameters: $\{\alpha, \beta\}$.

The third scheme was introduced in Teixeira (2001) for subtropical boundary layer clouds. Teixeira arrived at a diagnostic relationship for cloud cover by equating a cloud production term from detrainment and a cloud erosion term from turbulent mixing with the environment. We can express the Teixeira scheme as

$$C_{\text{Teixeira}} \stackrel{\text{def}}{=} \frac{Dq_c}{2q_s(1 - \widehat{\text{RH}})K} \left(-1 + \sqrt{1 + \frac{4q_s(1 - \widehat{\text{RH}})K}{Dq_c}} \right), \quad (5)$$

where $\widehat{\text{RH}} \stackrel{\text{def}}{=} \min\{\text{RH}, 1 - 10^{-9}\}$ bounds relative humidity to $1 - 10^{-9}$ to ensure reasonable asymptotics, $q_s = q_s(T, p)$ is the saturation specific humidity (Lohmann et al., 2016), and $\{D, K\}$ are the detrainment rate and the erosion coefficient, which are the two tuning parameters of the Teixeira scheme.

Besides those three traditional schemes, we additionally train the three neural networks (cell-, neighborhood-, and column-based NNs) from Grundner et al. (2022) on the DYAMOND data. These three NNs receive their inputs either from the same grid cell, the vertical neighborhood of the grid cell, or the entire grid column. Thus, they differ in the amount of vertical locality that is assumed for cloud cover parameterization. As the ‘undersampling step’ has to be done at a cell-based level, we omit it when pre-processing the training data for the column-based NN. Nevertheless, the column-based NN is evaluated on the same validation set as all other models.

Now that we have introduced three semi-empirical cloud cover schemes, which can be used as baselines, we are ready to derive a hierarchy of data-driven cloud cover schemes.

3.2 Developing Parsimonious Models via Sequential Feature Selection

Our goal is to develop parameterizations for cloud cover that are not only performant, but also simple and interpretable. Providing many, possibly correlated features to a model may needlessly increase its complexity and allow the model to learn spurious links between its inputs and outputs (Nowack et al., 2020), impeding both interpretability (Molnar, 2020) and generalizability (Brunton et al., 2016). Therefore, we instead seek parsimonious models, starting without any features and carefully selecting and adding features to a given type of model (e.g., a second-order polynomial) in a sequential manner. The chosen feature is always the one that maximizes the model’s performance on a sufficiently large subset of the training set (see also Sec 5.1.1). With this sequential feature selection (SFS) approach, we discourage the choice of correlated features and enforce sparsity by selecting a controlled number of features that already lead to the desired performance. Another benefit is that by studying the order of selected variables, optionally with the corresponding performance gains, we can gather intuition and physical knowledge about the task at hand. On the way, we will obtain an approximation of the best-performing set of features for a given number of features. There is however no guarantee of it truly being the best-performing feature set due to the greedy nature of the feature selection algorithm, which decreases its computational cost.

3.2.1 Linear Models and Polynomials

We allow first-order (i.e., linear models), second-order, and third-order polynomials. The set of possible terms grows from 25 (see Sec 2.1) to 325 for the second-order and to 2925 for the third-order polynomials. To circumvent memory issues for the third-order polynomials, we restrict the pool of possible features to combinations of the ten most important ones. In addition to these ten features, we also include air pressure to be technically able to later assign a sample to a cloud regime (overall reducing the total amount of possible terms from 2925 to 364). The choice of the ten features is informed by the SFS NNs (Sec 3.2.2), which are able to select informative features for nonlinear models. Thereafter we run SFS and train all linear models and polynomials using the *SequentialFeatureSelector* and *LinearRegression* methods of scikit-learn (Pedregosa et al., 2011) to obtain sequences of models with up to ten features.

3.2.2 Neural Networks

We train a sequence of SFS NNs with up to ten features using the “mlxtend” Python package (Raschka, 2018). We additionally train an NN with all 24 features for comparison purposes. As our regression task is similar in nature (including the vertical locality assumptions it makes for the features), we use the “Q1 NN” model architecture from Grundner et al. (2022) for all SFS NNs. “Q1 NN”’s architecture has three hidden layers with 64 units each; it uses batch normalization and its loss function includes L^1 and L^2 -regularization terms following hyperparameter optimization.

Due to the flexibility of NNs, when combining SFS with NNs, we obtain a sequence of features that is not bound to a particular model structure. In Sec 3.2.1 and 3.3, we therefore reuse the SFS NN feature rankings for other nonlinear models to restrict their set of possible features. The combination of SFS with NNs also yields a tentative upper bound on the accuracy one can achieve with N features: If we assume that i) SFS provides the best set of features for a given number of features N ; and ii) the NNs are able to outperform all other models given their features, one would not be able to outperform the SFS NNs with the same number of features. Even though the assumptions are only met approximately, we still receive helpful upper bounds on the performance of any model with N features.

3.3 Symbolic Regression Fits

To improve upon the analytical models of Sec 3.1 and 3.2.1 without compromising interpretability, we use recently-developed symbolic regression packages. We choose the PySR (Cranmer, 2020) and default GP-GOMEA (Virgolin et al., 2021) libraries, which are both based on genetic programming. GP-GOMEA is one of the best symbolic regression libraries according to SRBench, a symbolic regression benchmarking project that compared 14 contemporary symbolic regression methods (La Cava et al., 2021). PySR is a very flexible, efficient, well-documented, and well-maintained library. In PySR, we choose a large number of potential operators to enable a wide range of functions (see Appendix B for details). We also tried AIFeynman and found that its underlying assumption that one could learn from the NN gradient was problematic for less idealized data. Other promising packages from the SRBench competition, such as DSR/DSO and (Py)Operon, are left for future work. PySR and GP-GOMEA can only utilize a very limited number of features. Regardless of the number of features we provide, GP-GOMEA only uses 3–4, while PySR uses 5–6 features. For this reason, PySR also has a built-in tree-based feature selection method to reduce the number of features. Since the SFS NNs from Sec 3.2.2 already provide a sequence of features that can be used in general, nonlinear cases, we instead select the first five of these features to maximize comparability between models. The decision to run PySR with five features is also motivated by the good performance ($R^2 > 0.95$) of the corresponding SFS NN (see Sec 5.1.2). Each run of the PySR or GP-GOMEA algorithms adds new candidates to the list of final equations. From ≈ 600 of resulting equations, we select those that have a good skill ($R^2 > 0.9$), are interpretable, and satisfy most of the physical constraints that we define in the following section. As a final step, we refine the free parameters in the equation using the Nelder-Mead and BFGS optimizers (as in Sec 3.1).

4 Model Evaluation

4.1 Physical Constraints

To facilitate their use, we postulate that simple equations for cloud cover $\mathcal{C}(X)$ ought to satisfy certain physical constraints (Gentine et al., 2021; Kashinath et al., 2021): 1) The cloud cover output should be between 0 and 100%; 2) an absence of cloud condensates should imply an absence of clouds; 3-5) when relative humidity or the cloud water/ice mixing ratios increase (keeping all other features fixed), then cloud cover should not decrease; 6) cloud cover should not increase when temperature increases; 7) the function should be smooth on the entire domain. We can mathematically formalize these physical constraints (PC):

- 1) $\text{PC}_1: \mathcal{C}(X) \in [0, 100]\%$
- 2) $\text{PC}_2: (q_c, q_i) = 0 \Rightarrow \mathcal{C}(X) = 0$
- 3) $\text{PC}_3: \partial\mathcal{C}(X)/\partial\text{RH} \geq 0$
- 4) $\text{PC}_4: \partial\mathcal{C}(X)/\partial q_c \geq 0$

- 324 5) PC_5 : $\partial\mathcal{C}(X)/\partial q_i \geq 0$
 325 6) PC_6 : $\partial\mathcal{C}(X)/\partial T \leq 0$
 326 7) PC_7 : $\mathcal{C}(X)$ is a smooth function

327 While these physical constraints are intuitive, they will not be respected by data-driven
 328 cloud cover schemes if they are not satisfied in the data. In the DYAMOND data, the
 329 first physical constraint is always satisfied, and PC_2 is satisfied in 99.7% of all condensate-
 330 free samples. The remaining 0.3% are due to noise induced during coarse-graining. In
 331 order to check whether PC_3 - PC_6 are satisfied in our subset of the coarse-grained DYA-
 332 MOND data, we extract $\{q_c, q_i, \text{RH}, T\}$. We then separate the variable whose partial deriva-
 333 tive we are interested in. Bounded by the min/max-values of the remaining three vari-
 334 ables, we define a cube in this three-dimensional space, which we divide into N^3 equally-
 335 sized cubes. In this way, the three variables of the samples within the cubes become more
 336 similar with increasing N . If we now fit a linear function in a given cube with the sep-
 337 arated variable as the inputs and cloud cover as the output, then we can use the sign of
 338 the function’s slope to know whether the physical constraint is satisfied.

339 On one hand, the test is more expressive the smaller the cubes are, as the samples
 340 have more similar values for three of the four chosen variables and we can better approx-
 341 imate the partial derivative with respect to the separated variable. However, we only guar-
 342 antee similarity in three variables (omitting e.g., pressure). On the other hand, as the
 343 size of the cubes decreases, so does the number of samples contained in a cube, and noisy
 344 samples may skew the results. We therefore only consider the cubes that contain a suf-
 345 ficiently large number of samples (at least 10^4 out of the $2.9 \cdot 10^8$).

Table 1. The percentage of data cubes that fulfill a given physical constraint. Only the cubes with a sufficiently large amount of samples are taken into account. The last column shows the proportion of cubes (across all sizes we consider) in which the constraint is satisfied on average.

	(Maximum) Number of data cubes							Average (%)
	1	2 ³	3 ³	4 ³	5 ³	6 ³	7 ³	
PC_3	100	100	100	100	100	100	100	100
PC_4	100	100	83	90	73	78	71	77.5
PC_5	100	100	85	50	81	83	68	73.8
PC_6	100	50	100	67	72	89	75	77.7

346 We collect the results in Table 1, and find that the physical constraint PC_3 (with
 347 respect to RH) is always satisfied. The other constraints are satisfied in most (on aver-
 348 age 76%) of the cubes. Thus, from the data we can deduce that the final cloud cover scheme
 349 should satisfy PC_1 - PC_3 in all and PC_4 - PC_6 in most of the cases.

To enforce PC_1 , we always constrain the output to $[0, 100]$ before computing the MSE. With the exception of the linear and polynomial SFS models, we already ensure PC_1 during training. For PC_2 , we can define cloud cover to be 0 if the grid cell is condensate-free. We can combine PC_1 and PC_2 as

$$\mathcal{C}(X) = \begin{cases} 0, & \text{if } q_i + q_c = 0 \\ \max\{\min\{f(X), 100\}, 0\}, & \text{otherwise,} \end{cases} \quad (6)$$

350 and our goal is to learn the best fit for $f(X)$. In the case of the Xu-Randall and Teix-
 351 eira schemes, ensuring PC_2 is not necessary since they satisfy the constraint by design.

352

4.2 Performance Metrics

We use different metrics to train and validate the cloud cover schemes. We always train to minimize the mean squared error (MSE), which directly measures the average squared mismatch of the predictions $f(x_i)$ (usually set to be in $[0, 100]$) and the corresponding true (cloud cover) values y_i :

$$\text{MSE} \stackrel{\text{def}}{=} \frac{1}{N} \sum_{i=1}^N (\mathcal{C}(x_i) - y_i)^2. \quad (7)$$

The coefficient of determination R^2 -value takes the variance of the output $Y = \{y_i\}_{i=1}^N$ into account:

$$R^2 \stackrel{\text{def}}{=} 1 - \frac{\text{MSE}}{\text{Var}(Y)}. \quad (8)$$

To compare discrete univariate probability distributions P and Q , we use the Hellinger distance

$$H(P, Q) \stackrel{\text{def}}{=} \frac{1}{\sqrt{2}} \|\sqrt{P} - \sqrt{Q}\|_2. \quad (9)$$

353

As opposed to the Kullback-Leibler divergence, the Hellinger distance between two distributions is always symmetric and finite (in $[0, 1]$).

354

355

As our measure of complexity we use the number of (free/tunable/trainable) parameters of a model. A clear limitation of this complexity measure is that, e.g., the expression $f(x) = ax$ is considered as complex as $g(x) = \sin(\exp(ax))$. However, in this study, most of our models (i.e., the linear models, polynomials, and NNs) do not contain these types of nested operators. Instead, each additional parameter usually corresponds to an additional term in the equation. In the case of symbolic regression tools, operators are already taken into account (see Appendix B) during the selection process, and we find that the number of trainable parameters suffices to compare the complexity of our symbolic equations in their simplified forms. Finally, this complexity measure is one of the few that can be used for both analytical equations and NNs.

356

357

358

359

360

361

362

363

364

4.3 Cloud Regime-Based Evaluation

We define four cloud regimes based on air pressure p and the total cloud condensate q_t (cloud water plus cloud ice) mixing ratio:

365

366

367

368

1. Low air pressure, little condensate (cirrus-type cloud regime)
2. High air pressure, little condensate (cumulus-type cloud regime)
3. Low air pressure, substantial condensate (deep convective-type cloud regime)
4. High air pressure, substantial condensate (stratus-type cloud regime)

369

370

371

372

373

374

375

376

377

378

Pressure or condensate values that are above their medians (78 787 Pa and $1.62 \cdot 10^{-5}$ kg/kg) are considered to be large, while values below the median are considered small. Each regime has a similar amount of samples (between 35 and 60 million samples per regime). In this simplified data split, based on Rossow and Schiffer (1991), air pressure and total cloud condensate mixing ratio serve as proxies for cloud top pressure and cloud optical thickness. These regimes will help decompose model error to better understand the strengths and weaknesses of each model, discussed in the following section.

379

5 Results

380

5.1 Performance on the Storm-Resolving (DYAMOND) Training Set

381

382

383

In this section, we train the models we introduced in Sec 3 on the (coarse-grained) DYAMOND training data and compare their performance and complexity on the DYAMOND validation data. We start with the sequential feature selection's results.

384

5.1.1 Feature Ranking

385

386

387

388

389

390

We perform 10 SFS runs for each linear model, polynomial, and NN from Sec 3.3. Each run varies the random training subset, which consists of $\mathcal{O}(10^5)$ samples in the case of NNs and $\mathcal{O}(10^6)$ samples in the case of polynomials. We then average the rank of a selected feature and note it down in brackets. We omit the average rank if it is the same for each random subset. By \mathcal{P}_d , $d \in \{1, 2, 3\}$ we denote polynomials of degree d (e.g., \mathcal{P}_1 groups linear models). The sequences in which the features are selected are

$$\mathcal{P}_1: \text{RH} \rightarrow T \rightarrow \partial_z \text{RH} \rightarrow q_i[4.3] \rightarrow \partial_{zz} p[4.7] \rightarrow q_c \rightarrow U \rightarrow \partial_{zz} q_c \rightarrow \partial_z q_v \rightarrow z_g$$

$$\mathcal{P}_2: \text{RH} \rightarrow T \rightarrow q_c q_i \rightarrow \text{RH} \partial_z \text{RH} \rightarrow T \partial_z \text{RH}[5.6] \rightarrow q_v \text{RH}[6.4] \rightarrow \text{TRH}[7.4] \rightarrow \text{RH}^2[7.9] \rightarrow \partial_z q_v[9.2] \rightarrow U[10.1]$$

$$\mathcal{P}_3: \text{RH} \rightarrow T \rightarrow q_c q_i \rightarrow T^2 \text{RH}[4.4] \rightarrow \text{RH}^2[5.4] \rightarrow T^2[6.7] \rightarrow \text{RH} \partial_z \text{RH}[7.4] \rightarrow \partial_z \text{RH}[8.3] \rightarrow p^2 \partial_{zz} p[8.8] \rightarrow T \partial_z \text{RH}[9.4]$$

$$\text{NNs: RH} \rightarrow q_i \rightarrow q_c \rightarrow T[4.1] \rightarrow \partial_z \text{RH}[4.9] \rightarrow \partial_{zz} p[6.7] \rightarrow \partial_z p[8.1] \rightarrow \partial_{zz} \text{RH}[8.3] \rightarrow \partial_z T[10.0] \rightarrow p_s[10.1]$$

391

392

393

394

395

396

397

398

399

400

401

Regardless of the model, the selection algorithm chooses RH as the most informative feature for predicting cloud cover. This is consistent with, e.g., Walcek (1994), who considers RH to be the best single indicator of cloud cover in most of the troposphere. Considering that the cloud cover in the high-resolution data was only derived from the cloud condensate mixing ratio, the models' prioritization of RH is quite remarkable. From the feature sequences, we can also deduce that cloud cover depends on the mixing ratios of cloud condensates in a very nonlinear way: The polynomials choose $q_i q_c$ as their third feature and do not use any other terms containing q_i or q_c . The NNs choose q_i and q_c as their second and third features, and are able to express a nonlinear function of these two features. The linear model cannot fully exploit q_i and q_c and hence attaches less importance to them.

402

403

404

405

406

407

408

409

410

411

412

Since RH and T are chosen as the most informative features for the linear model, we can derive a notable linear dependence of cloud cover on these two features (the corresponding model being $f(\text{RH}, T) = 41.31\text{RH} - 15.54T + 44.63$). However, given the possibility, higher order terms of T and RH are chosen as additional predictors over, for instance, p or q_v . Finally, $\partial_z \text{RH}$ is an important recurrent feature for all models. Depending on the model, the coefficient associated with $\partial_z \text{RH}$ can be either negative or positive. If $\partial_z \text{RH} \neq 0$, one can assume some variation of cloud cover (i.e., cloud area fraction) vertically within the grid cell. Thus, $\partial_z \text{RH}$ is a meaningful proxy for the subgrid vertical variability of cloud area fraction. Since the effective cloud area fraction of the entire grid cell is related to the maximum cloud area fraction at a given height within the grid cell, this could explain the significance of $\partial_z \text{RH}$.

413

5.1.2 Balancing Performance and Complexity

414

415

416

417

418

419

420

421

422

423

In Fig 2, we depict all of our models in a performance \times complexity plane. We measure performance as the MSE on the validation (sub)set of the DYAMOND data and use the number of free parameters in the model as our complexity metric. We add the Pareto frontier, defined to pass through the best-performing models of a given complexity. The SFS sequences described above are used to train the SFS models of the corresponding type. The only exception is the swapped order of $\partial_z p$ and $\partial_{zz} p$ for the NNs, as we base the sequence shown in Fig 2 on a single SFS run. For the SFS NNs with 4-7 features, it was possible to reduce the number of layers and hidden units without significant performance degradation, which reduced the number of free parameters by about an order of magnitude and put them on the Pareto frontier.

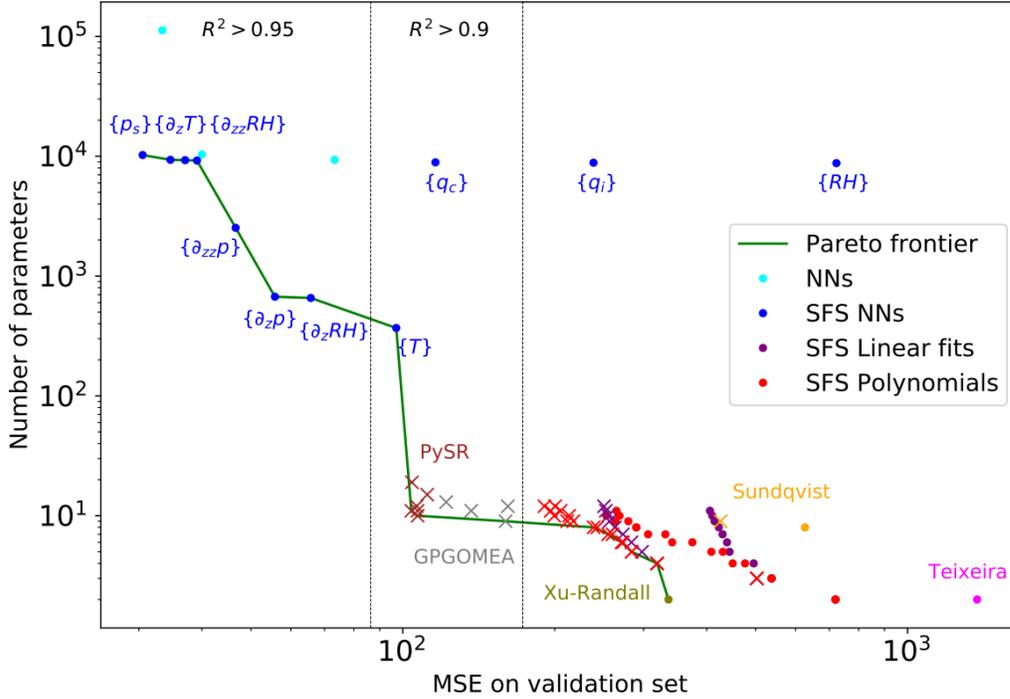


Figure 2. All models described in Sec 3 in a performance \times complexity plot. The dashed vertical lines mark the $R^2 = 0.95$ - and $R^2 = 0.9$ -boundaries. Models marked with a cross satisfy the second physical constraint PC_2 (using equation (6)). Only the best PySR and GP-GOMEA symbolic regression fits are shown. The NNs in cyan are the column-, neighborhood- and cell-based NNs when read from left to right. The SFS NN with the lowest MSE contains all 24 features described in Sec 2.1. For the SFS NNs, the last added feature is specified in curly brackets. Since the validation MSE of the SFS NNs decreases with additional features, we can extract the features for a given SFS NN by reading from right to left (e.g., the features of the SFS NN marked with $\{q_c\}$ are $\{q_i, q_c, RH\}$).

424 For most models, we train a second version that does not need to learn that condensate-
 425 free cells are always cloud-free, but for which the constraint is embedded by equation (6).
 426 For such models, condensate-free cells are removed from the training set. In addition to
 427 the schemes of Xu-Randall and Teixeira (see Sec 4.1), we find that it is also not neces-
 428 sary to enforce PC_2 in the case of NNs, since they are able to learn PC_2 without degrad-
 429 ing their performance. PC_1 is always enforced by default for all models.

430 We find that, even though the Sundqvist and Teixeira schemes are also tuned to
 431 the training set, linear models of the same complexity outperform them. However, these
 432 linear models do not lie on the Pareto frontier either. The lower performance of the Teix-
 433 eira scheme is most likely due to the fact that it was developed for subtropical bound-
 434 ary layer clouds. Among the existing schemes, only the Xu-Randall scheme with its two
 435 tuning parameters set to $\{\alpha, \beta\} = \{0.9, 9 \cdot 10^5\}$ is on the Pareto frontier as the sim-

436 plest model. With relatively large values for α and β , cloud cover is always approximately
 437 equal to relative humidity (i.e., $\mathcal{C} \approx \text{RH}^{0.9}$) when clouds are present. The next mod-
 438 els on the Pareto frontier are third-order SFS polynomials \mathcal{P}_3 with 2-6 features with PC_2
 439 enforced. To account for the bias term and the output of the polynomial being set to
 440 zero in condensate-free cells, the number of their parameters is the number of features
 441 plus 2. We then pass the line with $R^2 = 0.9$ and find three symbolic regression fits on
 442 the Pareto frontier, each trained on the five most informative features for the SFS NNs.
 443 All symbolic regression equations that appear in the plot are listed in Appendix C. We
 444 will analyze the PySR equation with arguably the best tradeoff between complexity (11
 445 free parameters when phrased in terms of normalized variables) and performance ($MSE =$
 446 $103.95(\%)^2$) in Sec 6. The remaining models on the Pareto frontier are SFS NNs with
 447 4-10 features and finally the NN with all 24 features defined in Sec 2.1 included ($MSE =$
 448 $30.51(\%)^2$).

449 Interestingly, the (quasi-local) 24-feature NN is able to achieve a slightly lower MSE
 450 ($30.51(\%)^2$) than the (non-local) column-based NN ($33.37(\%)^2$) with its 163 features.
 451 The two aspects that benefit the 24-feature NN are the additional information on the
 452 horizontal wind speed U and its derivatives, and the smaller number of condensate-free
 453 cells in its training set due to undersampling (Sec 2.1 and 3.1). The SFS NN with 10 fea-
 454 tures already shows very similar performance ($MSE = 34.64(\%)^2$) to the column-based
 455 NN with a (12 times) smaller complexity and fewer, more commonly accessible features.

456 Comparing the small improvements of the linear SFS models (up to $MSE = 250.43(\%)^2$)
 457 with the larger improvements of SFS polynomials (up to $MSE = 190.78(\%)^2$) with in-
 458 creasing complexity, it can be deduced that it is beneficial to include nonlinear terms in-
 459 stead of additional features in a linear model. For example, NNs require only three fea-
 460 tures to predict cloud cover reasonably well ($R^2 = 0.933$), and five features are suffi-
 461 cient to produce an excellent model ($R^2 = 0.962$) because they learn to nonlinearly trans-
 462 form these features.

463 The PySR equations can estimate cloud cover very well ($R^2 \in [0.935, 0.940]$). How-
 464 ever, while the PySR equations depend on five features, the NNs are able to outperform
 465 them with as few as four features ($R^2 = 0.944$). This suggests that the NNs learn bet-
 466 ter functional dependencies than PySR, as they do better with less information. How-
 467 ever, the improved performance of the NNs comes at the cost of additional complexity
 468 and greatly reduced interpretability.

469 5.2 Split by Cloud Regimes

470 In this section, we divide the DYAMOND data set into the four cloud regimes in-
 471 troduced in Sec 4.3. In Fig 3, we compare the cloud cover predictions of Pareto-optimal
 472 models (on Fig 2's Pareto frontier) with the actual cloud cover distribution in these regimes.
 473 We evaluate the models located at favorable positions on the Pareto frontier (at the be-
 474 ginning to maximize simplicity, at the end to maximize performance, or on some corners
 475 to optimally balance both). Of the two PySR equations, we consider the one with the
 476 lowest MSE (as in Sec 6 later).

477 In general, we find that the PySR equation (except in the cirrus regime) and the
 478 6-feature NN can reproduce the distributions quite well (Hellinger distances < 0.05),
 479 while the 24-feature NN shows excellent skill (Hellinger distances ≤ 0.015). However,
 480 all models have difficulty predicting the number of fully cloudy cells in all regimes (es-
 481 pecially in the regimes with fewer cloud condensates).

482 For the PySR equation (and also the 24-feature NN), the cirrus regime distribu-
 483 tion is the most difficult one to replicate. The Hellinger distances suggest that it is the
 484 model's functional form, and not its number of features that limits model performance
 485 in the cirrus regime. Indeed, the decrease in the Hellinger distance between the PySR

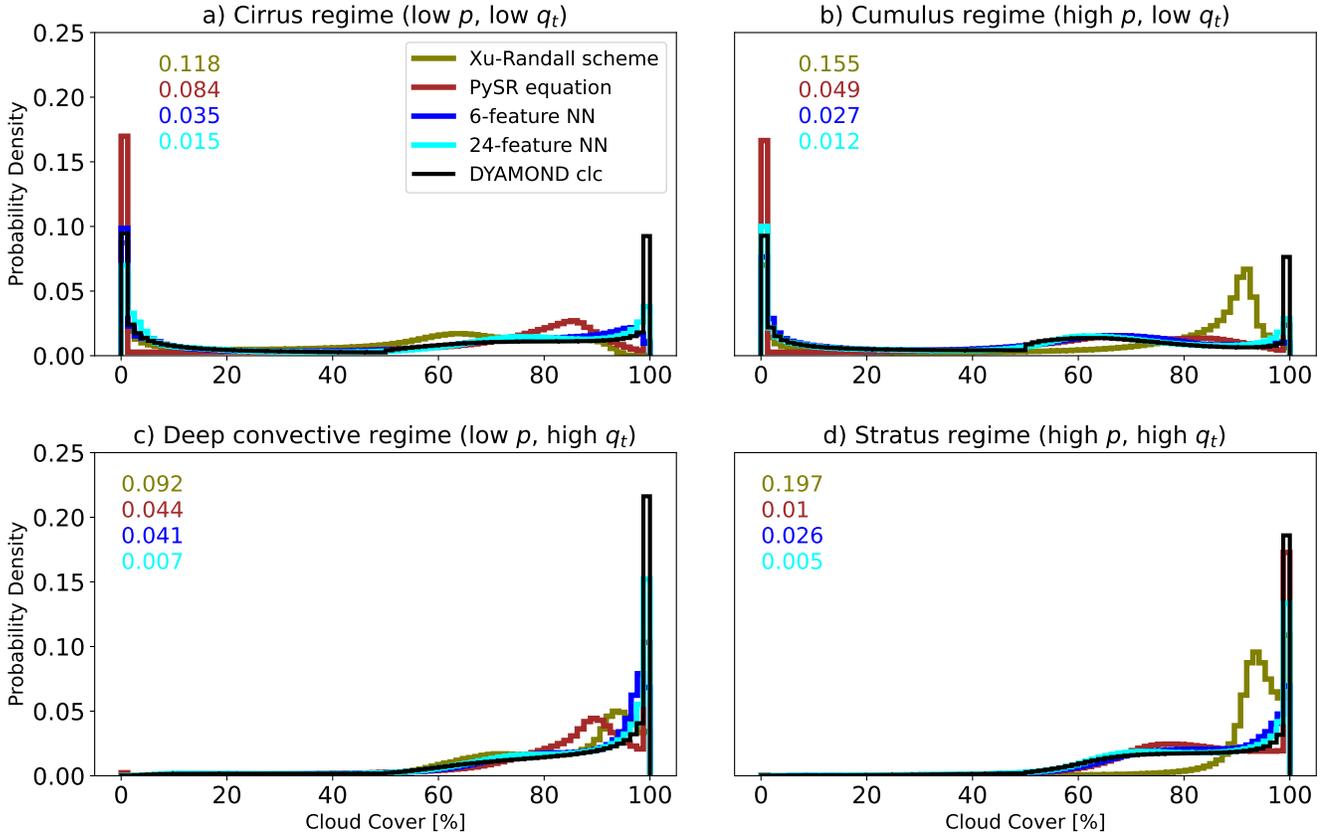


Figure 3. Predicted cloud cover distributions of selected Pareto-optimal models evaluated on the DYAMOND data, divided into four different cloud regimes. The numbers in the upper left indicate the Hellinger distance between the predicted and the actual cloud cover distributions for each model and cloud regime.

486 equation and the 6-feature NN is larger (0.049) than the decrease between the 6- and
 487 the 24-feature NN (0.02). Technically, the PySR equation has the same features as the
 488 5-feature and not the 6-feature NN, but the Hellinger distances of these two NNs to the
 489 actual cloud cover distribution are almost the same (difference of 0.003 in the cirrus regime).
 490 In the condensate-rich regimes, the PySR equation is as good as the 6-feature NN and
 491 even able to outperform it on the stratus regime. To improve the PySR scheme further
 492 in terms of its predicted cloud cover distributions, and combat its underestimation of cloud
 493 cover in the cirrus regime, it would most likely be a good strategy to sample more training
 494 data from the cirrus regime and less from the stratus regime. Note that the PySR
 495 equation actually achieves its best R^2 score ($R^2 = 0.84$) in the cirrus regime as the co-
 496 efficient of determination takes into account the high variance of cloud cover in the cir-
 497 rus regime.

498 Besides the peaks at the tail(s) of the distributions, we can see another erroneous
 499 peak in the Xu-Randall distribution in each cloud regime. We find that by neglecting
 500 the cloud condensate term and equating RH with the regime-based median, we can ap-
 501 proximately rederive these modes of the Xu-Randall cloud cover distributions in each
 502 regime using the Xu-Randall equation (4). With our choice of $\alpha = 0.9$, this mode is
 503 indeed very close (absolute difference at most 8% cloud cover) to the median relative hu-
 504 midity calculated in each regime. By increasing α , we should therefore be able to push
 505 the mode above 100% cloud cover and thus remove the spurious peak – however, at the
 506 cost of increasing the overall MSE of the Xu-Randall scheme.

507 5.3 Transferability to Meteorological Reanalysis (ERA5)

508 Designing data-driven models that are not specific to a given Earth system model
 509 and a given grid is challenging. Therefore, in this section, we aim to identify which of
 510 our Pareto-optimal ML models are most general and transferable. Furthermore, to our
 511 knowledge, there is no systematic method to incorporate observations into ML param-
 512 eterizations for climate modeling. We take a step towards transferring schemes trained
 513 on SRMs to observations by analyzing the ability of the Pareto-optimal schemes to trans-
 514 fer learn the ERA5 meteorological reanalysis from the DYAMOND set.

515 To do so, we take a certain number (either 1 or 100) of random locations, and col-
 516 lect the information from the corresponding grid columns of the ERA5 data over a cer-
 517 tain number of time steps in a data set \mathcal{T} . Starting from the parameters learned on the
 518 DYAMOND data, we retrain the cloud cover schemes on \mathcal{T} and evaluate them on the
 519 entire ERA5 data set. We can think of \mathcal{T} as mimicking a series of measurements at these
 520 locations, which help the schemes to adjust to the unseen data set. Fig 4 shows the MSE
 521 of the Pareto-optimal cloud cover schemes on the ERA5 data set after transfer learning
 522 on data sets \mathcal{T} of different sizes.

523 The first columns of the three panels show no variability because the schemes are
 524 applied directly to the ERA5 data without any transfer learning ($\mathcal{T} = \emptyset$). None of the
 525 schemes perform well without transfer learning ($R^2 < 0.15$), which is expected given
 526 the different distributions of cloud ice and water between the DYAMOND and ERA5
 527 data sets (Fig 1). That being said, the SFS NNs retain their superior performance (MSE
 528 $\approx 300(\%)^2$ without retraining), especially compared to the non-retrained SFS poly-
 529 nomials, which exhibit MSEs in the range of $1375 \pm 55(\%)^2$ and are therefore not shown
 530 in Panel c.

531 For most schemes, performance increases significantly after seeing one grid column
 532 of ERA5 data, with the exception of the SFS NNs with more than 6 features and the
 533 GPGOMEA equation. The performance of the GPGOMEA equation varies greatly be-
 534 tween the selected grid columns, and the SFS NNs with many features appear to under-
 535 fit the small transfer learning training set. The models with the lowest MSEs are (1) the
 536 slightly more complex of the two PySR equations (median MSE = $148(\%)^2$); and (2)

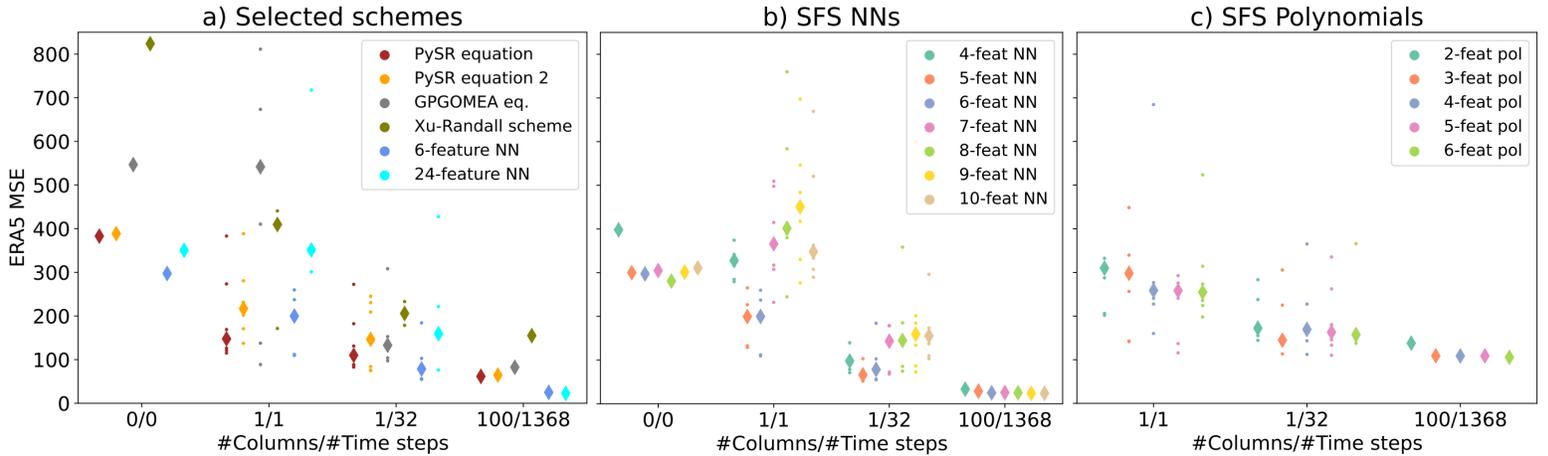


Figure 4. Performance of DYAMOND-trained Pareto-optimal cloud cover schemes on the ERA5 data set after transfer learning. The labels on the x-axis denote how many grid columns taken across how many time steps make up the transfer learning training set. Each setting is run with six different random seeds and the diamond-shaped markers indicate the respective medians.

537 the SFS NNs with 5 and 6 features (median $\text{MSE} = 200(\%)^2$). While we cannot con-
 538 firm that fewer features (5-6 features) help with off-the-shelf generalizability of the SFS
 539 NNs, they do improve the ability to transfer learn after seeing only a few samples from
 540 the ERA5 data.

541 After increasing the number of time steps to be included in \mathcal{T} to 32 (correspond-
 542 ing to one year of our preprocessed ERA5 data set), the performances of the models start
 543 to converge and the SFS NNs with 5 and 6 features and its large number of trainable
 544 parameters outperform the PySR equation (with median $\Delta\text{MSE} \approx 35(\%)^2$). From the
 545 last column we can conclude that a \mathcal{T} consisting of 100 columns from all available time
 546 steps is sufficient for the ERA5 MSE of all schemes to converge. Remarkably, the order
 547 from best- to worst-performing model is exactly the same as it was in Fig 2 on the DYA-
 548 MOND data set. Thus, we find that the ability to perform well on the DYAMOND data
 549 set is directly transferable to the ability to perform well on the ERA5 data set given enough
 550 data, despite fundamental differences between the data sets.

551 A useful property of a model is that it is able to transfer learn what it learned over
 552 an extensive initial dataset after tuning only on a few samples. We can quantify the abil-
 553 ity to transfer learn with few samples in two ways: First, we can directly measure the
 554 error on the entire data set after the model has seen only a small portion of the data (in
 555 our case the ERA5 MSEs of the 1/1-column). Second, if this error is already close to the
 556 minimum possible error of the model, then few samples are really enough for the model
 557 to transfer learn to the new data set (in our case, the difference of MSEs in the 1/1-column
 558 and the 100/1368-column). In terms of the first metric (MSEs in $(\%)^2$), the leading five
 559 models are the more complex PySR equation (147.6), the 5- and 6-feature NNs (199.6/199.8),
 560 the simpler PySR equation (216.8), and the 6-feature polynomial (254.6). In terms of
 561 the second metric (difference of MSEs in $(\%)^2$), the top five models are again the more
 562 complex PySR equation (86.0), the 6-, 5-, and 4-feature polynomials (149.1/149.4/150.5),

563 and the simpler PySR equation (152.3). If we add both metrics, weighing them equally,
 564 then the more complex PySR equation has the lowest inability to transfer learn with few
 565 samples (233.7), followed by the simpler PySR equation (369.1) and the 5- and 6-feature
 566 SFS NNs (370.5/374.5, where all numbers have units (%)²). As the more complex PySR
 567 equation is leading in both metrics, we can conclude that it is most able to transfer learn
 568 after seeing only one column of ERA5 data, and we further investigate its physical be-
 569 havior in the next section.

570 6 Physical Interpretation of the Best Analytical Scheme

We find that the two PySR equations on the Pareto frontier (see Fig 2) achieve a good compromise between accuracy and simplicity. Both satisfy most of the physical constraints that we defined in Sec 4.1. In this section, we analyze the (more complex) PySR equation with a lower validation MSE as we showed that it generalized best to ERA5 data (see Fig 4). We also conclude that the decrease in MSE is substantial enough ($\Delta\text{MSE} = 3.04\%^2$) to warrant the analysis of the (one parameter) more complex equation: The equation for the case with condensates can be phrased as

$$f(RH, T, \partial_z RH, q_c, q_i) = I_1(RH, T) + I_2(\partial_z RH) + I_3(q_c, q_i), \quad (10)$$

where

$$\begin{aligned} I_1(RH, T) &\stackrel{\text{def}}{=} a_1 RH^2 + (a_2 RH - a_3) T^2 - a_4 RHT + a_5 RH + a_6 T - a_7 \\ I_2(\partial_z RH) &\stackrel{\text{def}}{=} (a_8 \partial_z RH + a_9) (\partial_z RH)^2 \\ I_3(q_c, q_i) &\stackrel{\text{def}}{=} -1 / (a_{10} q_c + a_{11} q_i + \epsilon). \end{aligned}$$

To compute cloud cover in the general case, we plug equation (10) into equation (6), enforcing the first two physical constraints ($\mathcal{C}(X) \in [0, 100]\%$ and $\mathcal{C}(X) = 0$ in condensate-free cells). On the DYAMOND data we find the best values for the coefficients to be

$$\{a_1, \dots, a_{11}, \epsilon\} = \{203, 0.06588, 0.03969, 33.87, 4224.6, 18.9586, 2202.6, 2 \cdot 10^{10}, 6 \cdot 10^7, 8641, 32544, 0.0106\}.$$

The function $I_1(RH, T)$ is a quadratic polynomial of RH and T with $a_2 RHT^2$ as an additional cubic term. I_1 causes cloud cover to generally increase with relative humidity (Fig 5a and 6a). While I_1 does not ensure the constraint PC_6 ($\partial\mathcal{C}/\partial T \leq 0$) everywhere, cloud cover typically decreases with temperature for samples of the DYAMOND data set (see Fig 5f). In the hot limit, we have

$$\lim_{T \rightarrow \infty} I_1(RH, T) = \begin{cases} \infty, & RH > a_3/a_2 \\ -\infty, & RH < a_3/a_2. \end{cases}$$

Thus, in the case of DYAMOND, a relative humidity of $a_3/a_2 \approx 0.6$ defines a cutoff for cloudiness in the hot limit. To ensure PC_3 ($\partial\mathcal{C}/\partial RH \geq 0$) in all cases, we replace RH by

$$\max\{RH, -c_1 T^2 + c_2 T - c_3\}, \quad (11)$$

571 with $c_1 = a_2/(2a_1) \approx 0.00016$, $c_2 = a_4/(2a_1) \approx 0.0834$ and $c_3 = a_5/(2a_1) \approx 10.405$.
 572 We arrive at this expression when solving $\partial f/\partial RH$ for RH. This condition of replacing
 573 RH triggers in roughly 1% of our samples. It ensures that cloud cover does not increase
 574 when decreasing relative humidity in cases of low relative humidity and average temper-
 575 ature (see Fig 6). Modifying the equation (10) in such a way does not deteriorate its per-
 576 formance on the DYAMOND data. Fig 6b) illustrates how the modification ensures PC_3
 577 in an average setting. It would be difficult to apply a similar modification to the NN which
 578 in our case violates PC_3 for $RH > 0.95$. We can also identify another feature of equa-
 579 tion (10); the absence of a minimum value of relative humidity, below which cloud cover
 580 must always be zero (the *critical relative humidity threshold*).

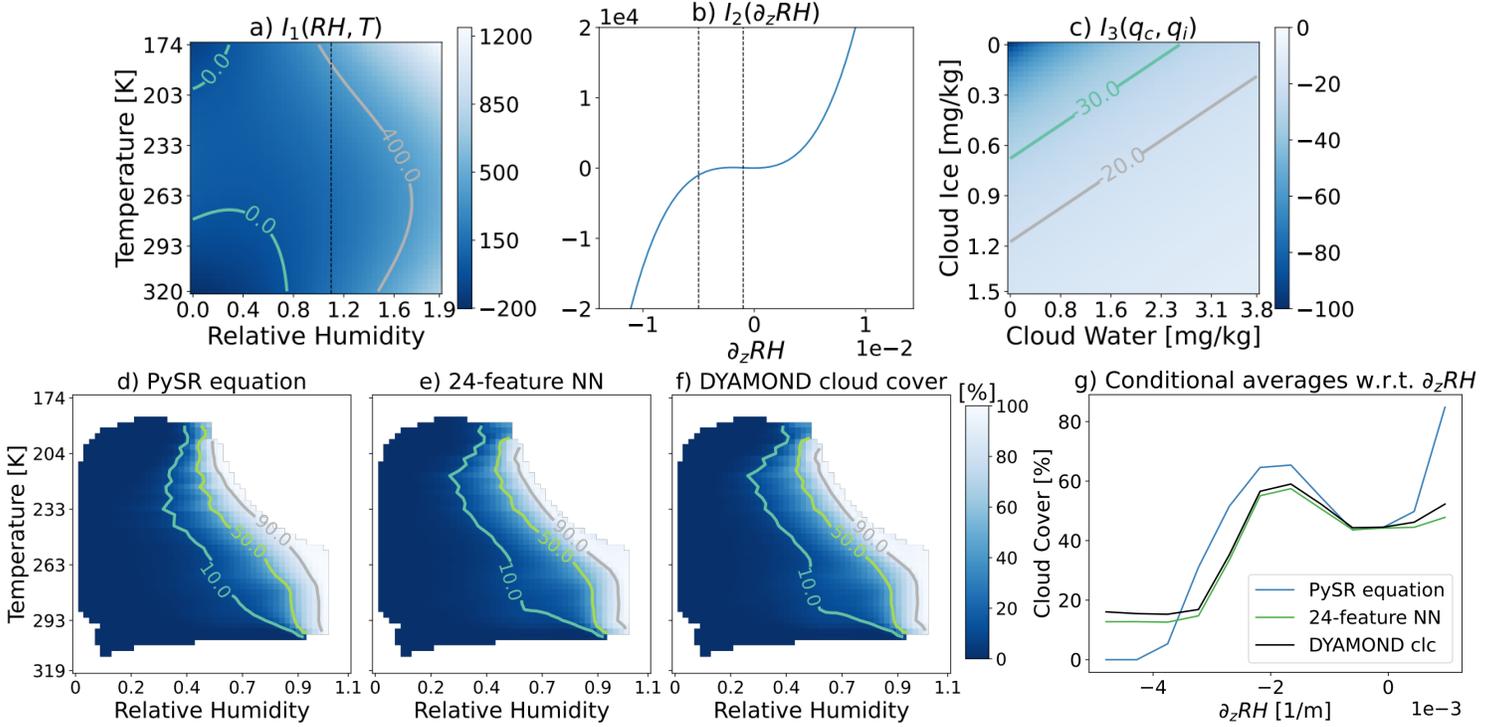


Figure 5. Top row: 1D- or 2D-plots of the three terms I_1, I_2, I_3 as functions of their inputs. In Panels a and b, the axis-values are bound by the respective minima and maxima in the DYAMOND data set, while those minima/maxima were divided by 5000 in Panel c. The vertical black lines indicate the region of values covered by Panels d-g. Bottom row: Conditional average plots of cloud cover with respect to relative humidity and temperature (Panels d-f) or $\partial_z RH$ (Panel g).

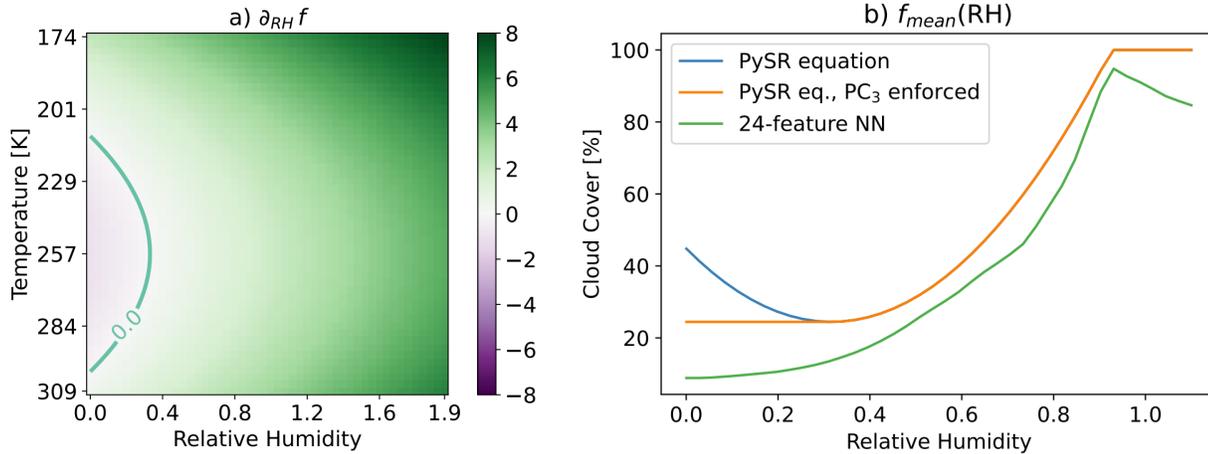


Figure 6. Panel a: Contour plot of $\partial_{RH}f$ as a function of relative humidity and temperature. The contour marks the boundary where $\partial_{RH}f = 0$. Panel b: Predictions of the PySR equation (10) with and without the modification (11) as a function of relative humidity. For comparison, the predictions of the SFS NN with 24 features are shown. The other features are set to their respective mean values.

581 The second function $I_2(\partial_z RH)$ is a cubic polynomial of $\partial_z RH$. Its slope is mainly
 582 driven by a_8 . The parameter a_9 , while being negligible for small values of $\partial_z RH$, increases
 583 the value of I_2 . The negative quotient of both parameters $-a_9/a_8$ defines the threshold
 584 for $\partial_z RH$, below which $I_2 < 0$.

585 Removing the I_2 -term, we find that the induced prediction error is largest, on av-
 586 erage, in situations that are i) relatively dry ($RH \approx 0.6$), ii) close to the surface ($z \approx$
 587 1000m), iii) over water (land fraction ≈ 0.1), iv) characterized by an inversion ($\partial_z T \approx$
 588 0.01), and v) have small values of $\partial_z RH$ ($\partial_z RH \approx -0.002$, which corresponds to the
 589 cloud cover peak in Fig 5g). Using our cloud regimes of Sec 5.2, we find the average ab-
 590 solute error is largest in the stratus regime (4% cloud cover). Indeed, by plotting the glob-
 591 ally averaged contributions of I_1 , I_2 and I_3 on a vertical layer at about 1500m altitude
 592 (Fig A1), we find that I_2 is most active in regions with low-level inversions where ma-
 593 rine stratocumulus clouds are abundant (Mauritsen et al., 2019). From this, we can infer
 594 that the SFS NN has chosen $\partial_z RH$ as a useful predictor to detect marine stratocu-
 595 muli and the symbolic regression algorithm has found a way to express this relationship
 596 mathematically. It is more informative than $\partial_z T$ (rank 10 in Sec 5.1.1), which would mea-
 597 sure the strength of an inversion more directly. The significance of $\partial_z RH$ relates back
 598 to our discussion in Sec 5.1.1: By comparing relative humidity values with those from
 599 the grid cells above and below we can infer cloud area fraction even when it is not rep-
 600 resented in the coarse variables of the given grid cell.

601 The third function $I_3(q_c, q_i)$ is always negative and decreases cloud cover where there
 602 is little cloud ice or water. It ensures that PC_4 and PC_5 are always satisfied. Large val-
 603 ues of a_{10} or a_{11} enable larger values of cloud water/ice to actually set I_3 close to zero.
 604 Finally, ϵ serves to avoid division by zero in condensate-free cells.

605 Given that equation (10) is a continuous function, the continuity constraint PC_7
 606 is only violated if and only if the cloud cover prediction is modified to be 0 in the condensate-
 607 free regime (by equation (6)), and would be positive otherwise. The value of ϵ dictates

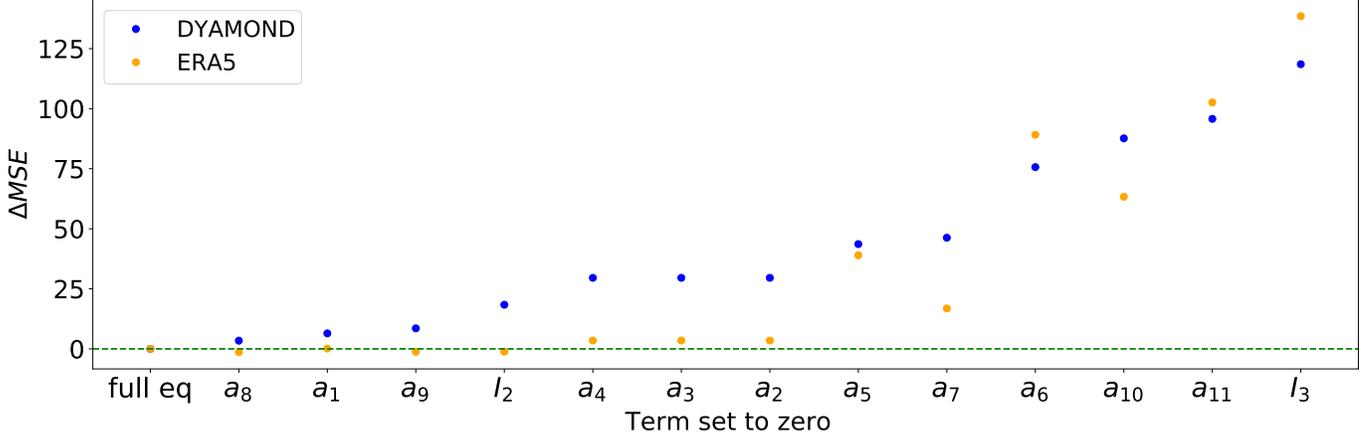


Figure 7. Ablation study of equation (10) on the DYAMOND and ERA5 data sets. The removal of the function I_1 leads to a large decrease of MSE (of $1300/763(\%)^2$) on the DYAMOND/ERA5 data sets and is therefore not shown.

608 how frequently the cloud cover prediction needs to be modified. In the limit $\epsilon \rightarrow 0$
 609 we could remove the different treatment of the condensate-free case. In our data set, equa-
 610 tion (10) yields a positive cloud cover prediction in 0.35% of condensate-free samples.
 611 Thus, the continuity constraint PC_7 is almost always satisfied (in 99.65% of our condensate-
 612 free samples).

613 To convince ourselves that all terms/parameters of equation (10) are indeed rel-
 614 evant to its skill, we examine the effects of their removal in an ablation study (Fig 7).
 615 We found that for the results to be meaningful, removing individual terms or param-
 616 eters requires readjusting the remaining parameters; in a setting with fixed param-
 617 eters the removal of multiple parameters often led to better outcomes than the removal of a
 618 single one of them. The optimizers (BFGS and Nelder-Mead) used to retune the remain-
 619 ing parameters show different success depending on whether the removal of terms is ap-
 620 plied to the equation formulated in terms of normalized or physical features (the latter
 621 being equation (10)). Therefore, each term is removed in both formulations, and the bet-
 622 ter result is chosen each time. To ensure robustness of the results, this ablation study
 623 is repeated for 10 different seeds on subsets with 10^6 data samples.

624 We find that the removal of any individual term in equation (10) would result in
 625 a noticeable reduction in performance on the DYAMOND data ($\Delta MSE \geq 3.4(\%)^2$ in
 626 absolute and $(MSE_{abl} - MSE_{full})/MSE_{abl} \geq 3.2\%$ in relative terms). Even though
 627 Fig 5g) suggests a cubic dependence of cloud cover on $\partial_z RH$, it is the least important
 628 term to include according to Fig 7. Applied to the ERA5 data, we can even dispense with
 629 the entire I_2 term. Furthermore, we find that the quadratic dependence on RH can be
 630 largely compensated by the linear terms. The most important terms to include are those
 631 with cloud ice/water and, concurring with the SFS polynomials in Sec 5.1.1, the linear
 632 dependence on RH and temperature. Coinciding with the SFS NN feature sequences in
 633 Sec 5.1.1, cloud ice ($\Delta MSE = 96/102(\%)^2$) is more important to take into account than
 634 cloud water ($\Delta MSE = 88/63(\%)^2$), especially for the ERA5 data set in which cloud
 635 ice is more abundant (see Fig 1). More generally, out of the functions I_1, I_2, I_3 we find
 636 $I_1(RH, T)$ to be most relevant ($\Delta MSE = 1300/763(\%)^2$), followed by $I_3(q_c, q_i)$ ($\Delta MSE =$
 637 $119/139(\%)^2$) and lastly $I_2(\partial_z RH)$ ($\Delta MSE = 18/-1(\%)^2$), once again matching the
 638 order of features that the SFS NNs had chosen.

7 Conclusion

In this study, we derived data-driven cloud cover parameterizations from coarse-grained global storm-resolving simulation (DYAMOND) output. We systematically populated a performance \times complexity plane with interpretable traditional parameterizations and regression fits on one side and high-performing neural networks on the other. Modern symbolic regression libraries (PySR, GPGOMEA) allow us to discover interpretable equations that diagnose cloud cover with excellent accuracy ($R^2 > 0.9$). From these equations, we propose a new analytical scheme for cloud cover (found with PySR) that balances accuracy ($R^2 = 0.94$) and simplicity (12 free parameters in the physical formulation). This analytical scheme satisfies six out of seven physical constraints (although the continuity constraint is violated in 0.35% of our condensate-free samples), providing the crucial third criterion for its selection. In a first evaluation, the (5-feature) analytical scheme was on par with the 6-feature NN in terms of reproducing cloud cover distributions (Hellinger distances < 0.05) in condensate-rich cloud regimes, yet underestimating cloud cover more strongly in condensate-poor regimes.

In addition to its interpretability, flexibility and efficiency, another major advantage of our best analytical scheme is its ability to adapt to a different data set (in our case, the ERA5 reanalysis product) after learning from only a few of the ERA5 samples in a transfer learning experiment. Due to the small amount of free parameters and the initial good fit on the DYAMOND data, our new analytical scheme outperformed all other Pareto-optimal models. We found that as the number of samples in the transfer learning sets increases, the models converged to the same performance rank on the ERA5 data as on the DYAMOND data, indicating strong similarities in the nature of the two data sets that could make which data set serves as the training set irrelevant. In an ablation study, we found that further reducing the number of free parameters in the analytical scheme would be inadvisable; all terms/parameters are relevant to its performance on the DYAMOND data. Key terms include a polynomial dependence on relative humidity and temperature, and a nonlinear dependence on cloud ice and water.

Our sequential feature selection approach with NNs revealed an objectively good subset of features for an unknown nonlinear function: relative humidity, cloud ice, cloud water, temperature and the vertical derivative of relative humidity (most likely linked to the vertical variability of cloud cover within a grid cell). While the first four features are well-known predictors for cloud cover, PySR also learned to incorporate $\partial_z RH$ in its equation. This additional dependence allows it to detect thin marine stratocumulus clouds, which are difficult, if not impossible to infer from exclusively local variables. These clouds are notoriously underestimated in the vertically coarse climate models (Nam et al., 2012). In ICON this issue is somewhat attenuated by multiplying, and thus increasing relative humidity in maritime regions by a factor depending on the strength of the low-level inversion (Mauritsen et al., 2019). Using symbolic regression, we thus found an alternative, arguably less crude approach, which could help mitigate this long-standing bias in an automated fashion. However, considering that this bias is still present in storm-resolving model simulations (Stevens et al., 2020), it could be advisable to further increase the resolution of the high-resolution model, and train on coarse-grained output from targeted large-eddy simulations (Stevens et al., 2005).

A crucial next step will be to test the cloud cover schemes when coupled to Earth system models, including ICON. We decided to leave this step for future work for several reasons. First, our focus was on the equation discovery methodology and the analysis of the discovered equation. Second, our goal was to derive a cloud cover scheme that is climate model-independent. Designing a scheme according to its online performance within a specific climate model decreases the likelihood of inter-model compatibility as the scheme has to compensate the climate model’s parameterizations’ individual biases. For instance, in ICON, the other parameterizations would most likely need to be re-calibrated to adjust for current compensating biases, such as clouds being ‘too few and too bright’

692 (Crueger et al., 2018). Third, the metrics used to validate a coupled model remain an
693 active research area, and at this point, it is unclear which targets must be met to accept
694 a new ML-based parameterization. That being said, the superior transferability of our
695 analytical scheme to the ERA5 reanalysis data not only suggests its applicability to ob-
696 servational data sets, but also that it may be transferable to other Earth system mod-
697 els.

698 Our current approach has some limitations. Symbolic regression libraries are lim-
699 ited in discovering equations with a large number of features. In many cases, five fea-
700 tures are insufficient to uncover a useful data-driven equation, requiring a reduction of
701 the feature space’s dimensionality. To measure model complexity, we used the number
702 of free parameters, disregarding the number of features and operators. Although the num-
703 ber of operators in our study was roughly equivalent to the number of parameters, this
704 may not hold in more general applications and the complexity of individual operators
705 would need to be specified (as in Appendix B).

706 Our approach differs from similar methods used to discover equations for ocean sub-
707 grid closures (Ross et al., 2023; Zanna & Bolton, 2020) because we included nonlinear
708 dependencies without assuming additive separability, instead fitting the entire equation
709 non-iteratively. Despite our efforts, the equation we found is still not as accurate as an
710 NN with equivalent features in the cirrus-like regime (the Hellinger distance between the
711 analytical scheme and the DYAMOND cloud cover distribution was more than twice as
712 large as for the NN). Comparing the partial dependence plots of the equation with those
713 of the NN could provide insights and define strategies to further extend and improve the
714 equation, while reducing the computational cost of the discovery. There are various meth-
715 ods available for utilizing NNs in symbolic regression for more than just feature selec-
716 tion, one of which is AIFeynman (Udrescu et al., 2020). While AIFeynman is based on
717 the questionable assumption that the gradient of an NN provides useful information, a
718 direct prediction of the equation using recurrent neural networks presents a promising
719 avenue for improved symbolic regression (Petersen et al., 2021; Tenachi et al., 2023).

720 Nonetheless, our simple cloud cover equation already achieves high performance.
721 Our study thus underscores that symbolic regression can complement deep learning by
722 deriving interpretable equations directly from data, suggesting untapped potential in other
723 areas of Earth system science and beyond.

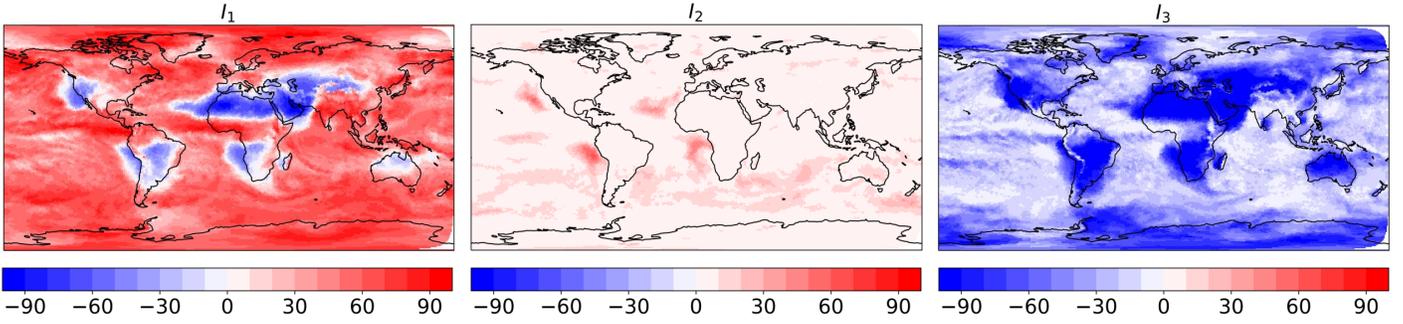


Figure A1. Maps of $I_1(RH, T)$, $I_2(\partial_z RH)$ and $I_3(q_c, q_i)$ on a vertical layer with an average height of 1490m, averaged over 10 days (11 Aug - 20 Aug, 2016). The data source is coarse-grained three-hourly DYAMOND data.

724 Appendix A Global Maps of I_1, I_2, I_3

725 In this section, we plot average function values for the three terms I_1, I_2 , and I_3
 726 of equation (10). We focus on the vertical layer roughly corresponding to an altitude of
 727 1500m to analyze if one of the terms would detect thin marine stratocumulus clouds. Due
 728 to their small vertical extent, these clouds are difficult to pick up on in coarse climate
 729 models, which constitutes a well-known bias. To compensate for this bias, the current
 730 cloud cover scheme of ICON has been modified so that relative humidity is artificially
 731 increased in low-level inversions over the ocean (Mauritsen et al., 2019).

732 Analyzing Fig A1, we find that the regions of high I_2 -values correspond with re-
 733 gions typical for low-level inversions and low-cloud fraction (Mauritsen et al., 2019; Muhl-
 734 bauer et al., 2014). These I_2 -values compensate partially negative I_1 - and I_3 -values in low-cloud
 735 regions of the Northeast Pacific, Southeast Pacific, Northeast Atlantic, and the South-
 736 east Atlantic. The I_1 -term is particularly small in the dry and hot regions of the Sahara
 737 and the Rub' al Khali desert and largest over the cold poles. The I_3 -term decreases cloud
 738 cover over land and is mostly inactive over the oceans due to the abundancy of cloud wa-
 739 ter.

740 Appendix B PySR Settings

First of all, we do not restrict the number of iterations, and instead restrict the run-
 time of the algorithm to ≈ 8 hours. We choose a large set of operators O to allow for
 various different functional forms (while leaving out non-continuous operators). To aid
 readability we show the operators applied to some $(x, y) \in \mathbb{R}^2$ which we denote by su-
 perscripts. To account for the different complexity of the operators, we split O into four
 distinct subsets

$$\begin{aligned}
 O_1^{(x,y)} &= \{x \cdot y, x + y, x - y, -x\} \\
 O_2^{(x,y)} &= \{x/y, |x|, \sqrt{x}, x^3, \max(0, x)\} \\
 O_3^{(x,y)} &= \{\exp(x), \ln(x), \sin(x), \cos(x), \tan(x), \sinh(x), \cosh(x), \tanh(x)\} \\
 O_4^{(x,y)} &= \{x^y, \Gamma(x), \operatorname{erf}(x), \arcsin(x), \arccos(x), \arctan(x), \operatorname{arsinh}(x), \operatorname{arcosh}(x), \operatorname{artanh}(x)\}
 \end{aligned}$$

741 of increasing complexity. The operators in $O_2/O_3/O_4$ are set to be 2/3/9 times as com-
 742 plex as those in O_1 . In this manner, for instance x^3 and $(x \cdot x) \cdot x$ have the same com-
 743 plexity. Furthermore, we assign a relatively low complexity to the operators in O_3 as they
 744 are very common and have well-behaved derivatives. With the factor of 9, we strongly
 745 discourage operators in O_4 . We expect that for every occurrence of a variable in a can-
 746 didate equation it will also need to be scaled by a certain factor. We do not want to dis-
 747 courage the use of such constant factors or the use of variables themselves and leave the
 748 complexity of constants and variables at their default complexity of one.

749 We obtain the best results when setting the complexity of the operators in O_1 to
 750 3 and training the PySR scheme on 5000 random samples. Other parameters include the
 751 population size (set to 20) and the maximum complexity of the equations that we ini-
 752 tially set to 200 and reduced to 90 in later runs.

753 Appendix C Selected Symbolic Regression Fits

This section lists all equations found with the symbolic regression libraries GP-GOMEA or PySR that are included in Fig 2, ranked in increasing MSE order. In brackets we provide the MSE/number of parameters. We list the equations according to their MSE. The equations that lie on the Pareto frontier are highlighted in bold:

1) PySR [103.95/11] :

$$f(\text{RH}, T, \partial_z \text{RH}, q_c, q_i) = \mathbf{203\text{RH}^2 + (0.06588\text{RH} - 0.03969)T^2 - 33.87\text{RHT} + 4224.6\text{RH} + 18.9586T - 2202.6 + (2 \cdot 10^{10} \partial_z \text{RH} + 6 \cdot 10^7)(\partial_z \text{RH})^2 - 1/(8641q_c + 32544q_i + 0.0106)}$$

2) PySR [104.26/19] :

$$f(\text{RH}, T, \partial_z \text{RH}, q_c, q_i) = (1.0364\text{RH} - 0.6782)(0.0581T - 16.1884)(-44639.6\partial_z \text{RH} + 1.1483T - 262.16) + 171.963\text{RH} - 1.4705T + 158.433(\text{RH} - 0.60251)^2 + (\partial_z \text{RH})^2(2 \cdot 10^{11}q_c - 8 \cdot 10^7\text{RH} + 7 \cdot 10^7) + 316.157 + 93319q_i - 1/(12108q_c + 39564q_i + 0.0111)$$

3) PySR [106.52/12] :

$$f(\text{RH}, T, \partial_z \text{RH}, q_c, q_i) = (57.2079\text{RH} - 34.4685)(3.0985\text{RH} + 73.1646(0.0039T - 1)^2 - 1.8669) + 123.175\text{RH} - 1.4091T + 1.5 \cdot 10^7(\partial_z \text{RH})^2(10619q_c - 4.9155\text{RH} + 4.7178) + 333.1 - 1/(10367q_c + 35939q_i + 0.0111)$$

4) PySR [106.95/11] :

$$f(\text{RH}, T, \partial_z \text{RH}, q_c, q_i) = 19.3885(3.0076\text{RH} - 1.8121)(3.2825\text{RH} + 73.1646(0.0039T - 1)^2 - 1.9777) + 118.59\text{RH} - 1.423T + 1.5 \cdot 10^7(3.0125 - 1.0129\text{RH})(\partial_z \text{RH})^2 + 339.2 - 1/(9325q_c + 34335q_i + 0.0109)$$

5) PySR [106.99/10] :

$$f(\text{RH}, T, \partial_z \text{RH}, q_c, q_i) = \mathbf{(58.189\text{RH} - 35.0596)(3.3481\text{RH} + 73.1646(0.0039T - 1)^2 - 2.0172) + 116.873\text{RH} - 1.4211T + 3.6 \cdot 10^7(\partial_z \text{RH})^2 + 339.9 - 1/(9237q_c + 34136q_i + 0.0109)}$$

6) PySR [111.76/15] :

$$f(\text{RH}, T, \partial_z \text{RH}, q_c, q_i) = (3.2665\text{RH} - 2.9617)(0.0435T - 9.0274)(16073.2\partial_z \text{RH} + 0.3013T - 68.4342) + 97.5754\text{RH} - 0.6556T + 175 + 123823q_i - 1/(9853q_c + 36782q_i + 0.0112)$$

7) GP-GOMEA [121.89/13] :

$$f(\text{RH}, T, q_c, q_i) = 8.459 \exp(2.559\text{RH}) - 33.222 \sin(0.038T + 109.878) + 24.184 - \sin(3.767\sqrt{|98709q_i - 0.334|})/(30046q_i + 5628q_c + 0.01)$$

8) GP-GOMEA [136.64/11] :

$$f(\text{RH}, T, q_c, q_i) = (8.65\text{RH} - 0.22T - 93.14)\sqrt{|0.62T - 414.23|} + 2368 - 1/(28661q_i + 4837q_c + 0.01)$$

9) GP-GOMEA [159.80/9] :

$$f(\text{RH}, q_c, q_i) = \mathbf{0.009e^{8.725\text{RH}} + 12.795 \log(229004q_i + 0.774(e^{11357q_c} - 1)) - 178246q_c + 66}$$

10) GP-GOMEA [161.45/12] :

$$f(\text{RH}, T, q_c, q_i) = (0.028e^{6.253\text{RH}} + 5\text{RH} - 0.076T + 4)/(183894q_i + 0.73e^{6565q_c - 91207q_i} - 0.62) + 92.3$$

754 Note that the assessed number of parameters is based on a simplified form of the
 755 equations in terms of its normalized variables. The amount of parameters in a given equa-
 756 tion is at least equal to the assessed number of parameters minus one (accounting for
 757 the zero in the condensate-free setting).

758 Data Availability Statement

759 The cloud cover schemes and analysis code can be found at [https://github.com/EyringMLClimateGroup/
 760 grundner23james_EquationDiscovery.CloudCover](https://github.com/EyringMLClimateGroup/grundner23james_EquationDiscovery.CloudCover) and are preserved at DOI:10.5281/
 761 zenodo.7817392. DYAMOND data management was provided by the German Climate
 762 Computing Center (DKRZ) and supported through the projects ESiWACE and ESiWACE2.
 763 The coarse-grained model output used to train and evaluate the neural networks amounts
 764 to several TB and can be reconstructed with the scripts provided in the GitHub repos-
 765 itory. The software code for the ICON model is available from [https://code.mpimet
 766 .mpg.de/projects/iconpublic](https://code.mpimet.mpg.de/projects/iconpublic).

767 Acknowledgments

768 Funding for this study was provided by the European Research Council (ERC) Synergy
 769 Grant “Understanding and Modelling the Earth System with Machine Learning (USMILE)”
 770 under the Horizon 2020 research and innovation programme (Grant agreement No. 855187).
 771 Beucler acknowledges funding from the Columbia University sub-award 1 (PG010560-
 772 01). Gentine acknowledges funding from the NSF Science and Technology Center, Cen-
 773 ter for Learning the Earth with Artificial Intelligence and Physics (LEAP) (Award 2019625).
 774 This manuscript contains modified Copernicus Climate Change Service Information (2023)
 775 with the following datasets being retrieved from the Climate Data Store: ERA5, ERA5.1
 776 (neither the European Commission nor ECMWF is responsible for any use that may be
 777 made of the Copernicus Information or Data it contains). The projects ESiWACE and
 778 ESiWACE2 have received funding from the European Union’s Horizon 2020 research and
 779 innovation programme under grant agreements No 675191 and 823988. This work used
 780 resources of the Deutsches Klimarechenzentrum (DKRZ) granted by its Scientific Steer-
 781 ing Committee (WLA) under project IDs bk1040, bb1153 and bd1179.

782 References

- 783 Beucler, T. G., Ebert-Uphoff, I., Rasp, S., Pritchard, M., & Gentine, P. (2022). Ma-
 784 chine learning for clouds and climate. *Earth Space Sci. Open Arch.*
- 785 Brenowitz, N. D., & Bretherton, C. S. (2018). Prognostic validation of a neural net-
 786 work unified physics parameterization. *Geophysical Research Letters*, *45*(12),
 787 6289-6298. doi: 10.1029/2018gl078510
- 788 Brunton, S. L., Proctor, J. L., & Kutz, J. N. (2016). Discovering governing equa-
 789 tions from data by sparse identification of nonlinear dynamical systems. *Pro-
 790 ceedings of the national academy of sciences*, *113*(15), 3932–3937.
- 791 Champion, K., Lusch, B., Kutz, J. N., & Brunton, S. L. (2019). Data-driven dis-
 792 covery of coordinates and governing equations. *Proceedings of the National
 793 Academy of Sciences*, *116*(45), 22445–22451.
- 794 Cranmer, M. (2020, September). *Pysr: Fast & parallelized symbolic regression in
 795 python/julia*. Zenodo. Retrieved from [http://doi.org/10.5281/zenodo
 796 .4041459](http://doi.org/10.5281/zenodo.4041459) doi: 10.5281/zenodo.4041459
- 797 Crueger, T., Giorgetta, M. A., Brokopf, R., Esch, M., Fiedler, S., Hohenegger, C.,
 798 ... others (2018). Icon-a, the atmosphere component of the icon earth system
 799 model: Ii. model evaluation. *Journal of Advances in Modeling Earth Systems*,
 800 *10*(7), 1638–1662.
- 801 Duras, J., Ziemann, F., & Klocke, D. (2021). The dyamond winter data collection. In
 802 *Egu general assembly conference abstracts* (pp. EGU21–4687).
- 803 Eyring, V., Mishra, V., Griffith, G. P., Chen, L., Keenan, T., Turetsky, M. R.,

- 804 ... van der Linden, S. (2021). Reflections and projections on a decade
805 of climate science. *Nature Climate Change*, *11*(4), 279-285. doi: 10.1038/
806 s41558-021-01020-x
- 807 Gao, F., & Han, L. (2012). Implementing the nelder-mead simplex algorithm with
808 adaptive parameters. *Computational Optimization and Applications*, *51*(1),
809 259-277.
- 810 Gentine, P., Eyring, V., & Beucler, T. (2021). Deep learning for the parametrization
811 of subgrid processes in climate models. *Deep Learning for the Earth Sciences:
812 A Comprehensive Approach to Remote Sensing, Climate Science, and Geo-
813 sciences*, 307-314.
- 814 Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., & Yacalis, G. (2018). Could
815 machine learning break the convection parameterization deadlock? *Geophysical
816 Research Letters*, *45*(11), 5742-5751.
- 817 Giorgetta, M. A., Crueger, T., Brokopf, R., Esch, M., Fiedler, S., Hohenegger, C.,
818 ... Stevens, B. (2018). Icon-a, the atmosphere component of the icon earth
819 system model: I. model description. *Journal of Advances in Modeling Earth
820 Systems*, *10*(7), 1638-1662. doi: 10.1029/2017ms001233
- 821 Giorgetta, M. A., Sawyer, W., Lapillonne, X., Adamidis, P., Alexeev, D., Clément,
822 V., ... Stevens, B. (2022). The icon-a model for direct qbo simulations on
823 gpus (version icon-cscs:baf28a514). *Geoscientific Model Development*, *15*(18),
824 6985-7016. Retrieved from [https://gmd.copernicus.org/articles/15/
825 6985/2022/](https://gmd.copernicus.org/articles/15/6985/2022/) doi: 10.5194/gmd-15-6985-2022
- 826 Grundner, A., Beucler, T., Gentine, P., Iglesias-Suarez, F., Giorgetta, M. A., &
827 Eyring, V. (2022). Deep learning based cloud cover parameterization for icon.
828 *Journal of Advances in Modeling Earth Systems*, *14*(12), e2021MS002959. doi:
829 <https://doi.org/10.1029/2021MS002959>
- 830 Hersbach, H., Bell, B., Berrisford, P., Biavati, G., Horányi, A., Muñoz Sabater, J.,
831 ... others (2018). Era5 hourly data on pressure levels from 1979 to present.
832 *Copernicus climate change service (c3s) climate data store (cds)*. (accessed at
833 DKRZ on 02-01-2023) doi: 10.24381/cds.bd0915c6
- 834 Hohenegger, C., Kornblueh, L., Klocke, D., Becker, T., Cioni, G., Engels, J. F., ...
835 Stevens, B. (2020). Climate statistics in global simulations of the atmosphere,
836 from 80 to 2.5 km grid spacing. *Journal of the Meteorological Society of Japan*,
837 *98*(1), 73-91. doi: 10.2151/jmsj.2020-005
- 838 Kashinath, K., Mustafa, M., Albert, A., Wu, J., Jiang, C., Esmailzadeh, S., ...
839 others (2021). Physics-informed machine learning: case studies for weather
840 and climate modelling. *Philosophical Transactions of the Royal Society A*,
841 *379*(2194), 20200093.
- 842 Krasnopolsky, V. M., Fox-Rabinovitz, M. S., & Belochitski, A. A. (2013). Using
843 ensemble of neural networks to learn stochastic convection parameterizations
844 for climate and numerical weather prediction models from data simulated by a
845 cloud resolving model. *Advances in Artificial Neural Systems*, *2013*, 1-13. doi:
846 10.1155/2013/485913
- 847 Kumar, I. E., Venkatasubramanian, S., Scheidegger, C., & Friedler, S. (2020). Prob-
848 lems with shapley-value-based explanations as feature importance measures. In
849 *International conference on machine learning* (pp. 5491-5500).
- 850 La Cava, W., Orzechowski, P., Burlacu, B., de Franca, F., Virgolin, M., Jin, Y., ...
851 Moore, J. (2021). Contemporary symbolic regression methods and their relative
852 performance. In J. Vanschoren & S. Yeung (Eds.), *Proceedings of the neu-
853 ral information processing systems track on datasets and benchmarks* (Vol. 1).
854 Retrieved from [https://datasets-benchmarks-proceedings.neurips.cc/
855 paper/2021/file/c0c7c76d30bd3dcaefc96f40275bdc0a-Paper-round1.pdf](https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/file/c0c7c76d30bd3dcaefc96f40275bdc0a-Paper-round1.pdf)
- 856 Lohmann, U., Lüönd, F., & Mahrt, F. (2016). *An introduction to clouds: From the
857 microscale to climate*. Cambridge University Press.
- 858 Lohmann, U., & Roeckner, E. (1996). Design and performance of a new cloud mi-

- 859 crophysics scheme developed for the echam general circulation model. *Climate*
860 *Dynamics*. doi: <https://doi.org/10.1007/BF00207939>
- 861 Mauritsen, T., Bader, J., Becker, T., Behrens, J., Bittner, M., Brokopf, R., ...
862 Roeckner, E. (2019). Developments in the mpi-m earth system model ver-
863 sion 1.2 (mpi-esm1.2) and its response to increasing co 2. *Journal of Advances*
864 *in Modeling Earth Systems*, 11(4), 998-1038. doi: 10.1029/2018ms001400
- 865 McCandless, T., Gagne, D. J., Kosović, B., Haupt, S. E., Yang, B., Becker, C., &
866 Schreck, J. (2022). Machine learning for improving surface-layer-flux estimates.
867 *Boundary-Layer Meteorology*, 185(2), 199–228.
- 868 Molnar, C. (2020). *Interpretable machine learning*. Lulu.com.
- 869 Molnar, C., Casalicchio, G., & Bischl, B. (2021). Interpretable machine learning—a
870 brief history, state-of-the-art and challenges..
- 871 Muhlbauer, A., McCoy, I. L., & Wood, R. (2014). Climatology of stratocumulus
872 cloud morphologies: microphysical properties and radiative effects. *Atmo-*
873 *spheric Chemistry and Physics*, 14(13), 6695–6716.
- 874 Nam, C., Bony, S., Dufresne, J.-L., & Chepfer, H. (2012). The ‘too few, too
875 bright’ tropical low-cloud problem in cmip5 models. *Geophysical Research*
876 *Letters*, 39(21).
- 877 Nocedal, J., & Wright, S. J. (1999). *Numerical optimization*. Springer.
- 878 Nowack, P., Runge, J., Eyring, V., & Haigh, J. D. (2020). Causal networks for cli-
879 mate model evaluation and constrained projections. *Nature communications*,
880 11(1), 1–11.
- 881 O’Gorman, P. A., & Dwyer, J. G. (2018). Using machine learning to parameter-
882 ize moist convection: Potential for modeling of climate, climate change, and
883 extreme events. *Journal of Advances in Modeling Earth Systems*, 10(10),
884 2548-2563. doi: 10.1029/2018ms001351
- 885 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ...
886 others (2011). Scikit-learn: Machine learning in python. *Journal of machine*
887 *learning research*, 12(Oct), 2825–2830.
- 888 Petersen, B. K., Landajuela, M., Mundhenk, T. N., Santiago, C. P., Kim, S. K., &
889 Kim, J. T. (2021). Deep symbolic regression: Recovering mathematical expres-
890 sions from data via risk-seeking policy gradients. In *Proc. of the international*
891 *conference on learning representations*.
- 892 Pincus, R., & Stevens, B. (2013). Paths to accuracy for radiation parameterizations
893 in atmospheric models. *Journal of Advances in Modeling Earth Systems*, 5(2),
894 225-233. doi: 10.1002/jame.20027
- 895 Raschka, S. (2018). Mlxtend: Providing machine learning and data science utilities
896 and extensions to python’s scientific computing stack. *The Journal of Open*
897 *Source Software*, 3(24). Retrieved from [http://joss.theoj.org/papers/](http://joss.theoj.org/papers/10.21105/joss.00638)
898 [10.21105/joss.00638](http://joss.theoj.org/papers/10.21105/joss.00638) doi: 10.21105/joss.00638
- 899 Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid
900 processes in climate models. *Proceedings of the National Academy of Sciences*,
901 115(39), 9684–9689.
- 902 Ross, A., Li, Z., Perezhogin, P., Fernandez-Granda, C., & Zanna, L. (2023).
903 Benchmarking of machine learning ocean subgrid parameterizations in an
904 idealized model. *Journal of Advances in Modeling Earth Systems*, 15(1),
905 e2022MS003258.
- 906 Rossow, W. B., & Schiffer, R. A. (1991). Isccp cloud data products. *Bulletin of the*
907 *American Meteorological Society*, 72(1), 2–20.
- 908 Rudy, S. H., Brunton, S. L., Proctor, J. L., & Kutz, J. N. (2017). Data-driven dis-
909 covery of partial differential equations. *Science advances*, 3(4), e1602614.
- 910 Schmidt, M., & Lipson, H. (2009). Distilling free-form natural laws from experimen-
911 tal data. *science*, 324(5923), 81–85.
- 912 Schulzweida, U. (2019, October). *Cdo user guide*. doi: 10.5281/zenodo.3539275

- 913 Smits, G. F., & Kotanchek, M. (2005). Pareto-front exploitation in symbolic regression.
914 *Genetic programming theory and practice II*, 283–299.
- 915 Stensrud, D. J. (2009). *Parameterization schemes: Keys to understanding numerical*
916 *weather prediction models*. Cambridge University Press.
- 917 Stevens, B., Acquistapace, C., Hansen, A., Heinze, R., Klinger, C., Klocke, D., ...
918 others (2020). The added value of large-eddy and storm-resolving models for
919 simulating clouds and precipitation. *Journal of the Meteorological Society of*
920 *Japan. Ser. II*, 98(2), 395–435.
- 921 Stevens, B., Moeng, C.-H., Ackerman, A. S., Bretherton, C. S., Chlond, A., de
922 Roode, S., ... others (2005). Evaluation of large-eddy simulations via obser-
923 vations of nocturnal marine stratocumulus. *Monthly weather review*, 133(6),
924 1443–1462.
- 925 Stevens, B., Satoh, M., Auger, L., Biercamp, J., Bretherton, C. S., Chen, X., ...
926 others (2019). Dyamond: the dynamics of the atmospheric general circula-
927 tion modeled on non-hydrostatic domains. *Progress in Earth and Planetary*
928 *Science*, 6(1), 1–17.
- 929 Sundqvist, H., Berge, E., & Kristjánsson, J. E. (1989). Condensation and cloud
930 parameterization studies with a mesoscale numerical weather prediction model.
931 *Monthly Weather Review*.
- 932 Teixeira, J. (2001). Cloud fraction and relative humidity in a prognostic cloud frac-
933 tion scheme. *Monthly Weather Review*, 129(7), 1750–1753.
- 934 Tenachi, W., Ibata, R., & Diakogiannis, F. I. (2023). Deep symbolic regression
935 for physics guided by units constraints: toward the automated discovery of
936 physical laws. *arXiv preprint arXiv:2303.03192*.
- 937 Trenberth, K. E., Fasullo, J. T., & Kiehl, J. (2009). Earth’s global energy budget.
938 *Bulletin of the American Meteorological Society*, 90(3), 311–324.
- 939 Udrescu, S.-M., Tan, A., Feng, J., Neto, O., Wu, T., & Tegmark, M. (2020). Ai feyn-
940 man 2.0: Pareto-optimal symbolic regression exploiting graph modularity. *Ad-*
941 *vances in Neural Information Processing Systems*, 33, 4860–4871.
- 942 Virgolin, M., Alderliesten, T., Witteveen, C., & Bosman, P. A. N. (2021). Improving
943 model-based genetic programming for symbolic regression of small expressions.
944 *Evolutionary Computation*, 29(2), 211–237.
- 945 Walcek, C. J. (1994). Cloud cover and its relationship to relative humidity during a
946 springtime midlatitude cyclone. *Monthly weather review*, 122(6), 1021–1035.
- 947 Wang, Y., Yang, S., Chen, G., Bao, Q., & Li, J. (2023). Evaluating two diagnos-
948 tic schemes of cloud-fraction parameterization using the cloudsat data. *Atmo-*
949 *spheric Research*, 282, 106510.
- 950 Weisman, M. L., Skamarock, W. C., & Klemp, J. B. (1997). The resolution de-
951 pendence of explicitly modeled convective systems. *Monthly Weather Review*,
952 125(4), 527–548.
- 953 Xu, K.-M., & Randall, D. A. (1996). A semiempirical cloudiness parameterization
954 for use in climate models. *Journal of the atmospheric sciences*, 53(21), 3084–
955 3102.
- 956 Zanna, L., & Bolton, T. (2020). Data-driven equation discovery of ocean mesoscale
957 closures. *Geophysical Research Letters*, 47(17), e2020GL088376.
- 958 Zhang, S., & Lin, G. (2018). Robust data-driven discovery of governing physical
959 laws with error bars. *Proceedings of the Royal Society A: Mathematical, Physi-*
960 *cal and Engineering Sciences*, 474(2217), 20180305.

Figure 1.

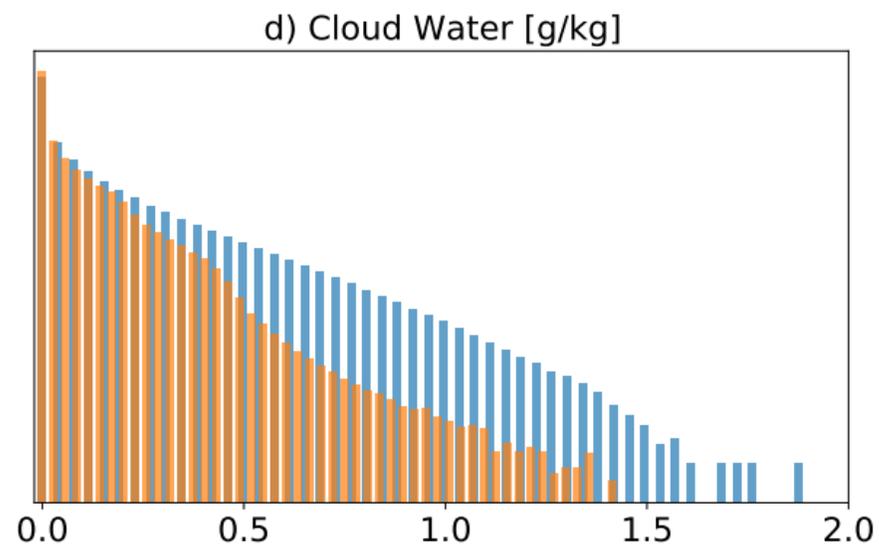
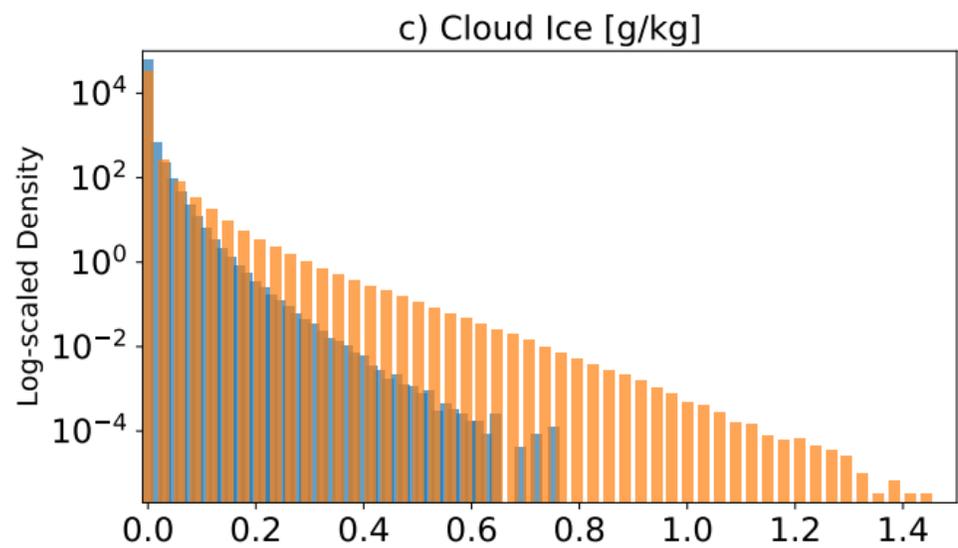
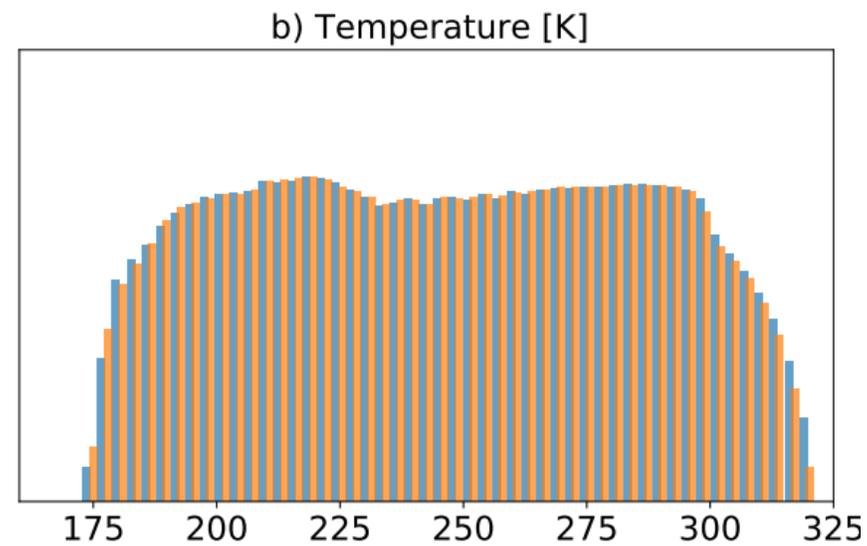
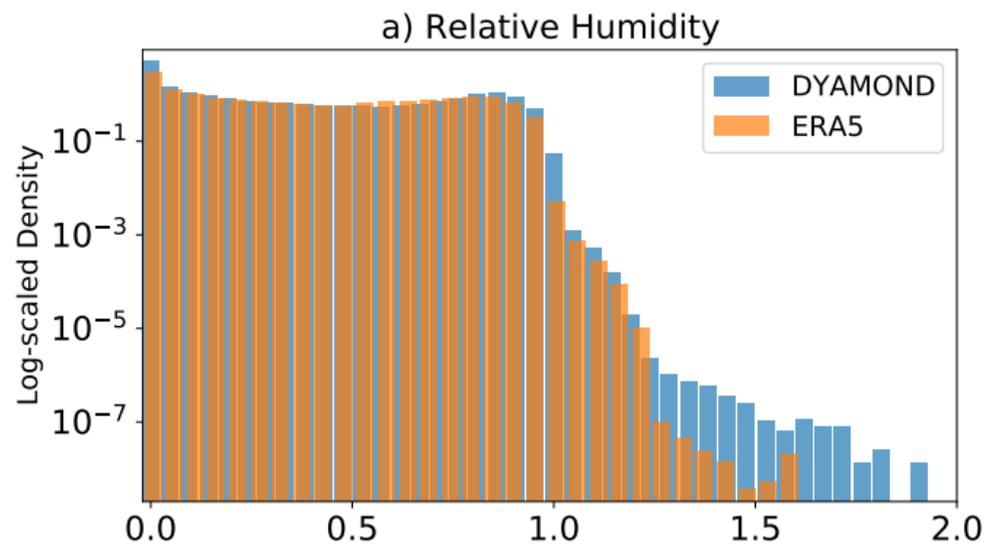


Figure 2.

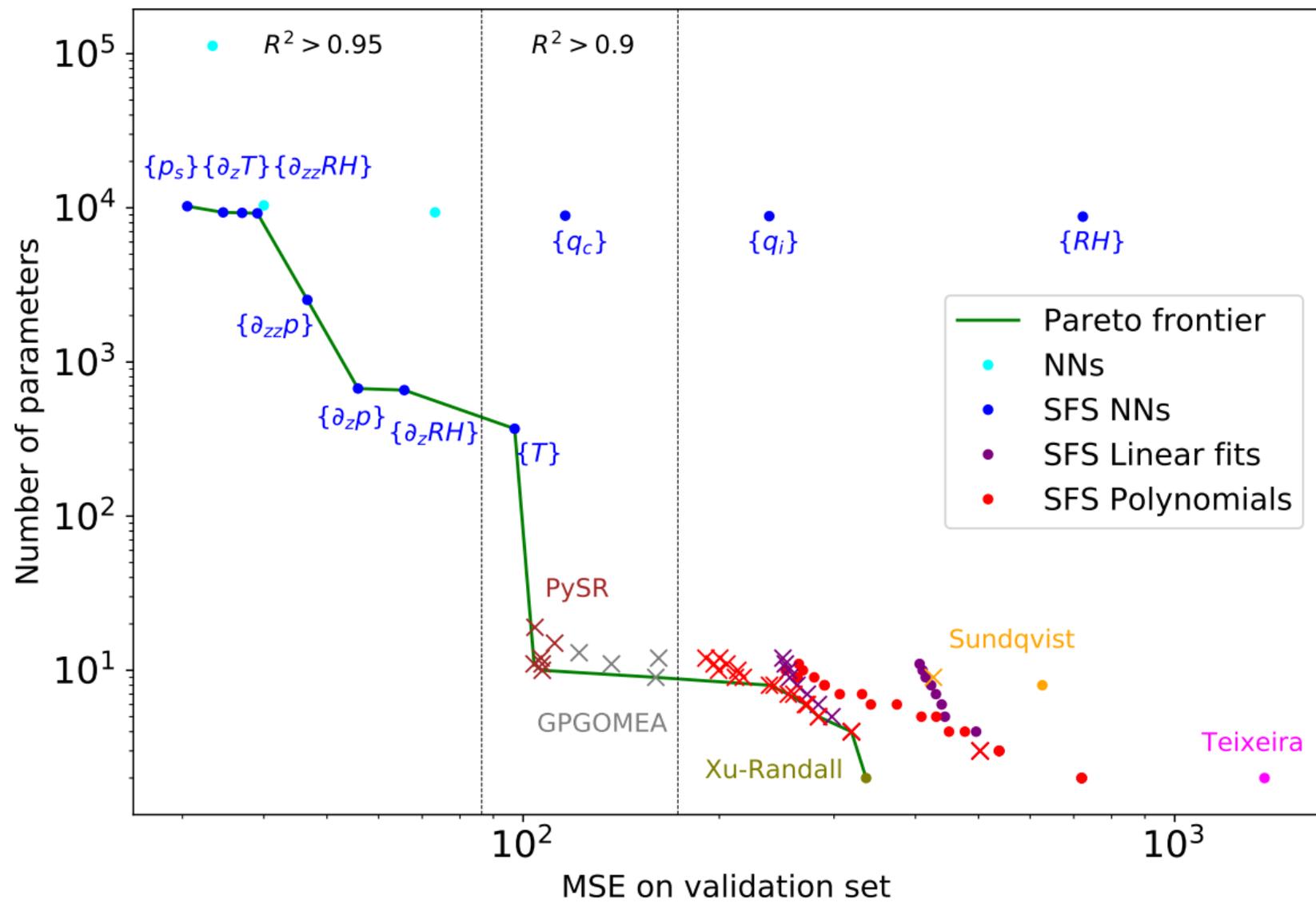
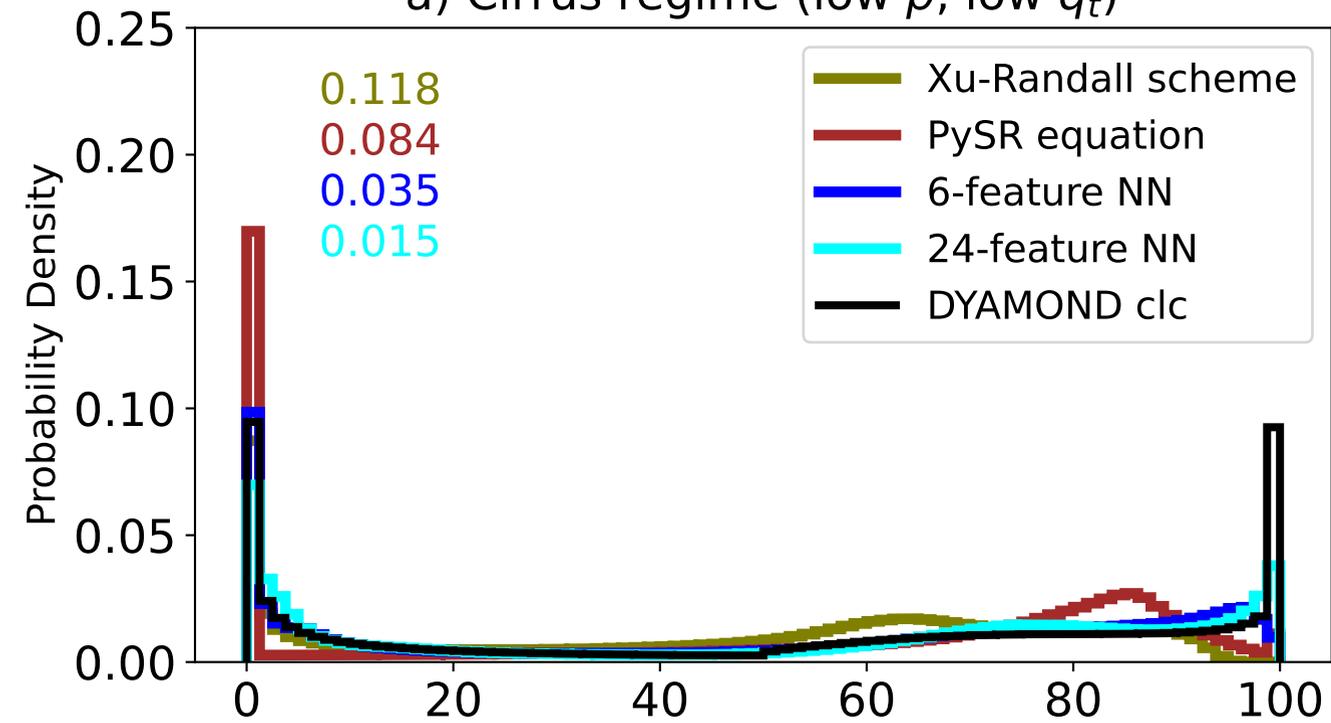
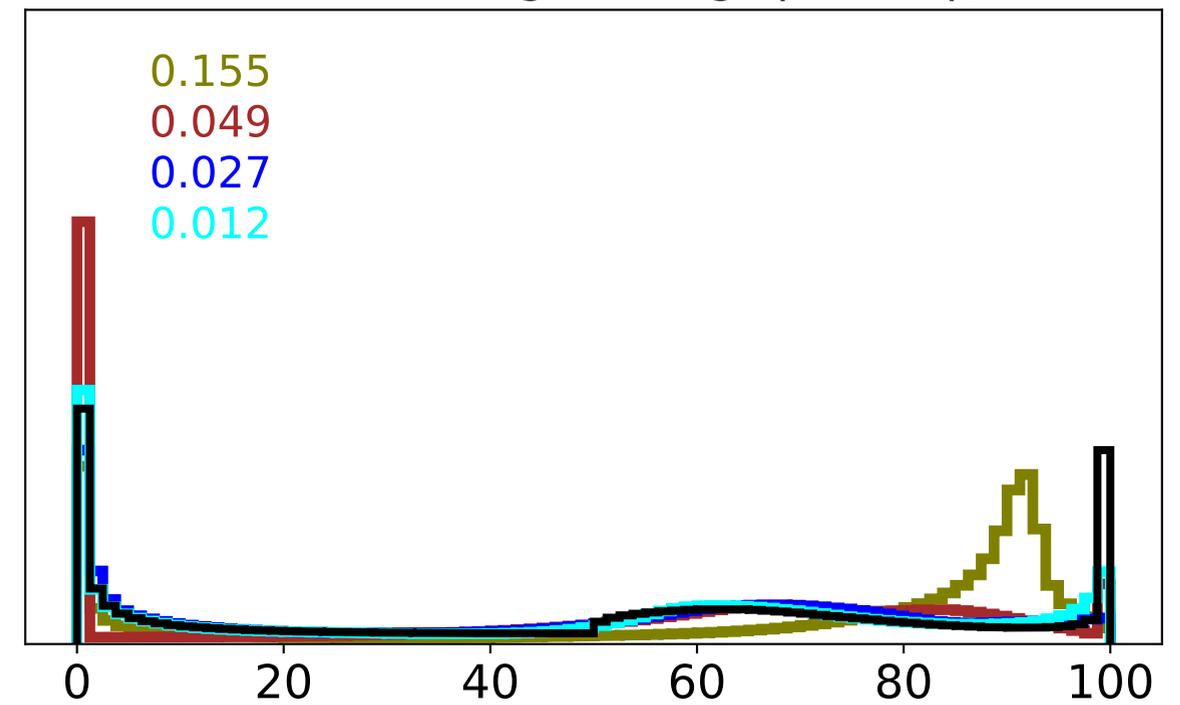


Figure 3.

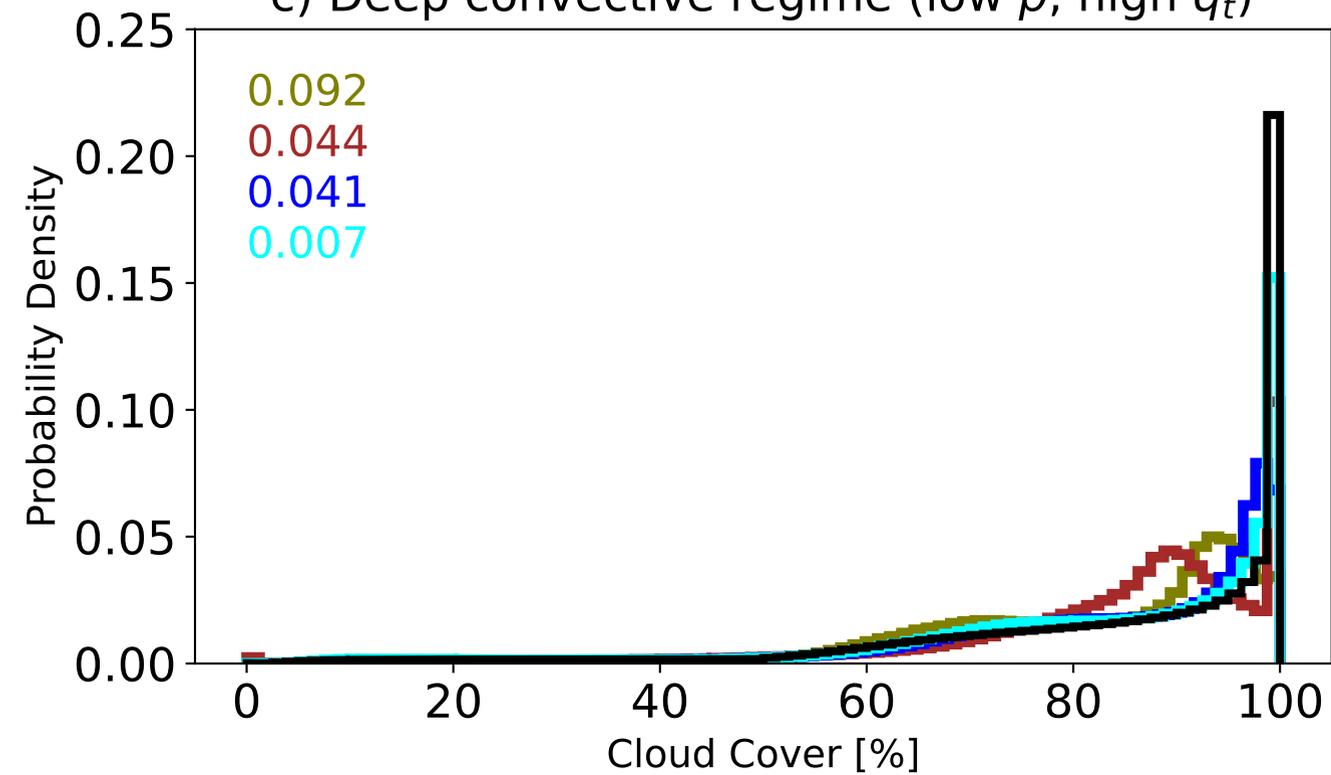
a) Cirrus regime (low p , low q_t)



b) Cumulus regime (high p , low q_t)



c) Deep convective regime (low p , high q_t)



d) Stratus regime (high p , high q_t)

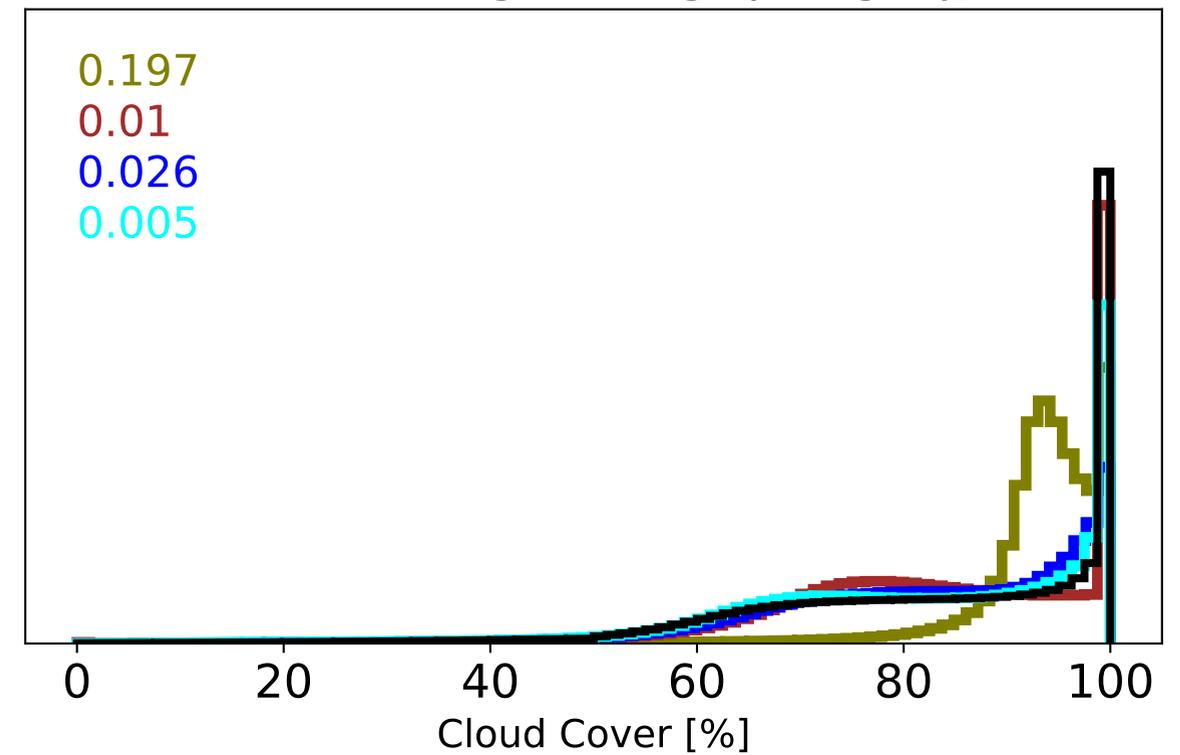
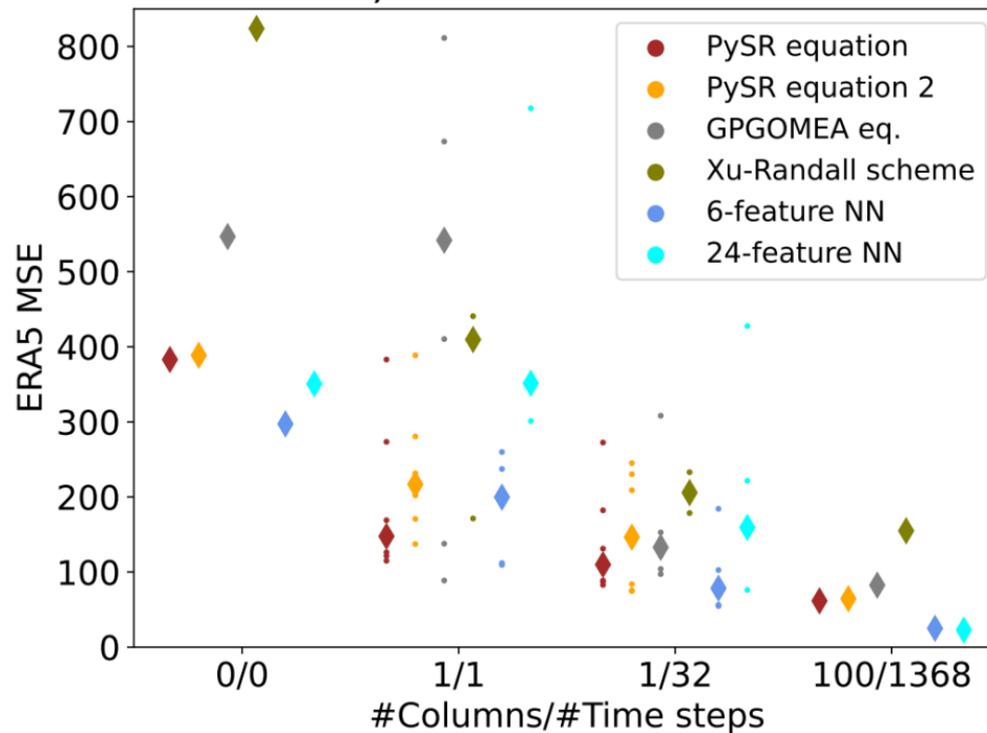
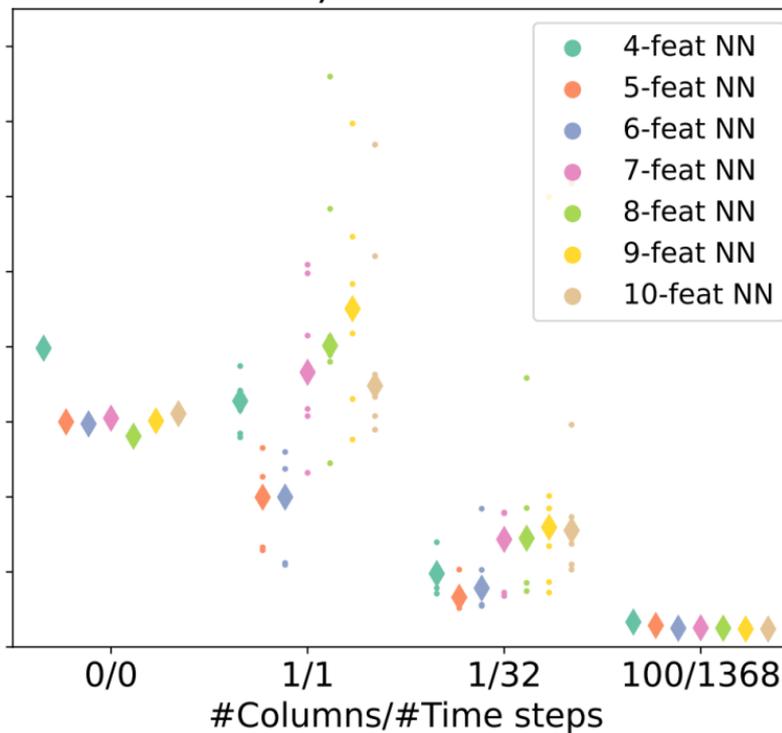


Figure 4.

a) Selected schemes



b) SFS NNs



c) SFS Polynomials

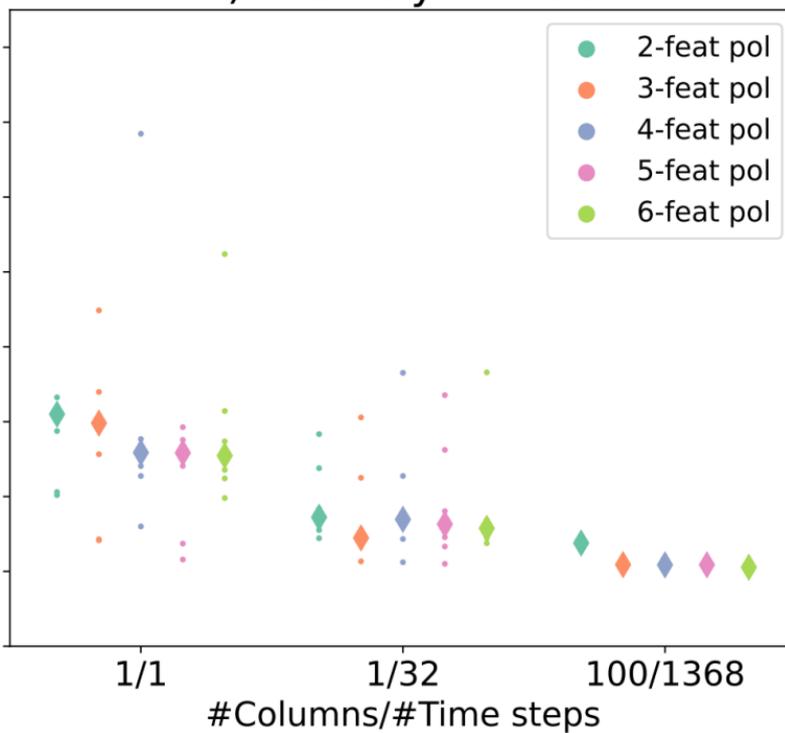


Figure 5.

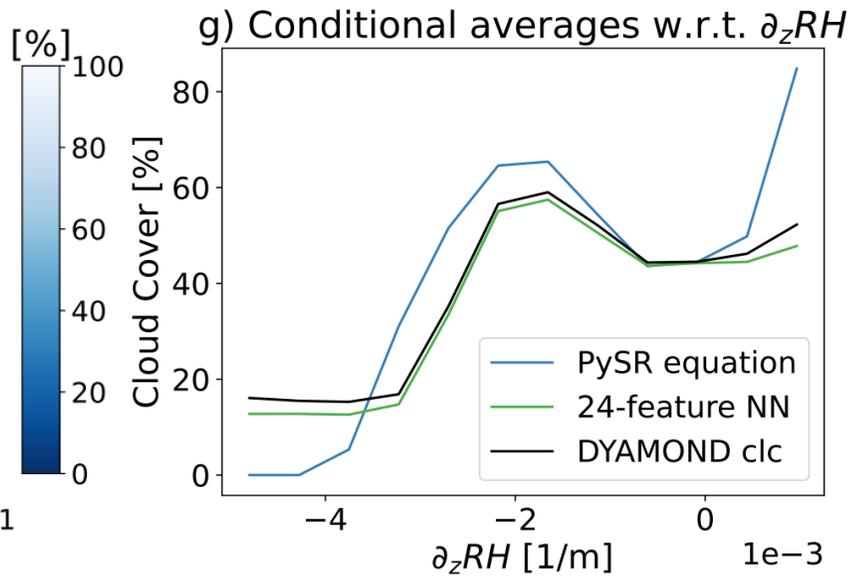
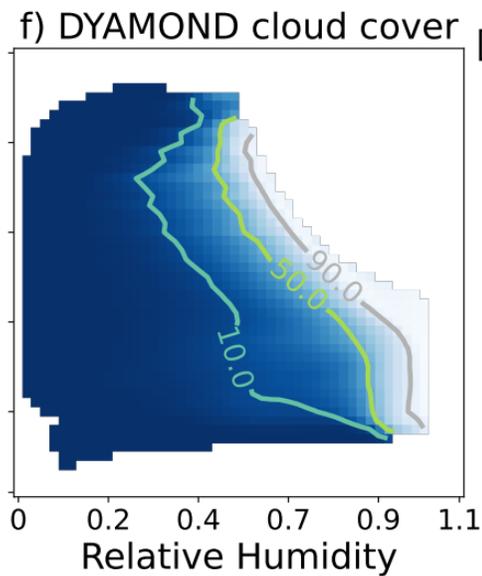
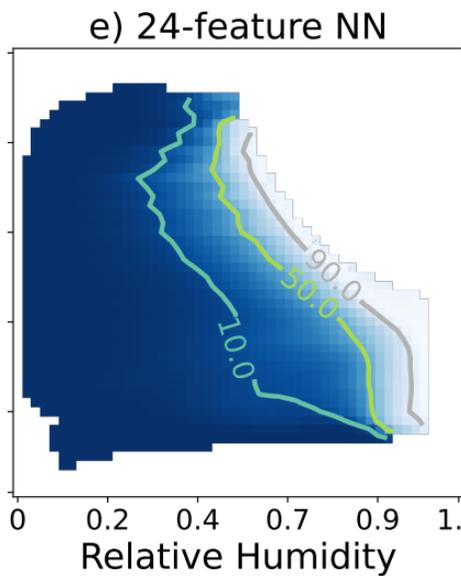
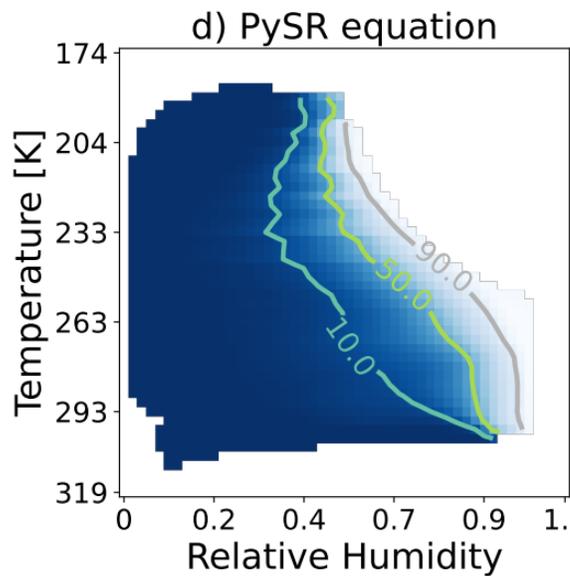
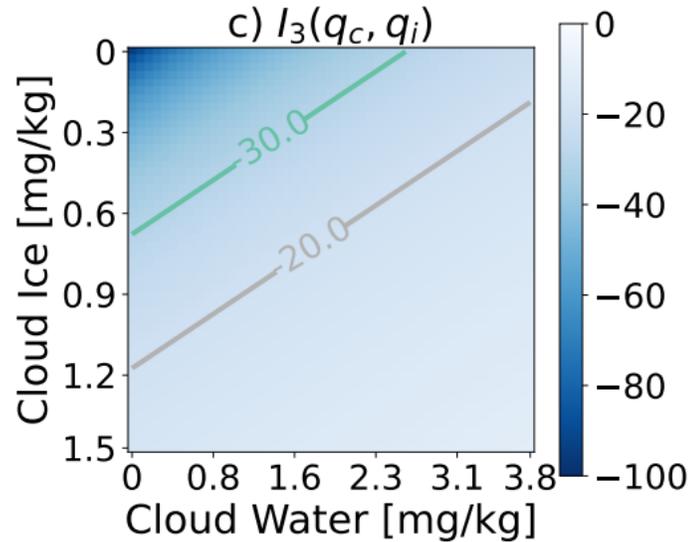
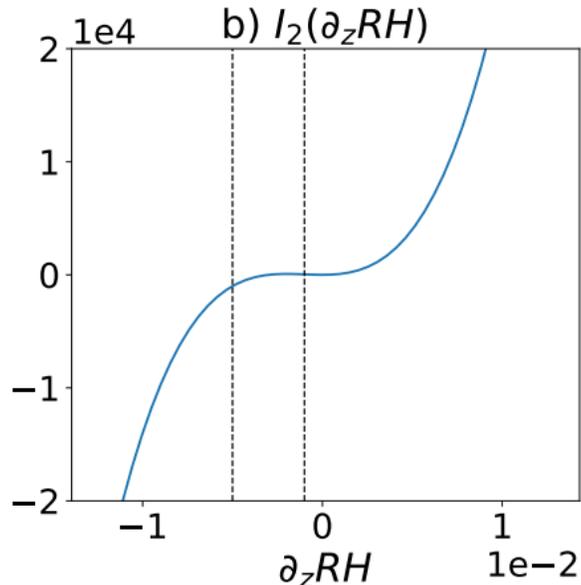
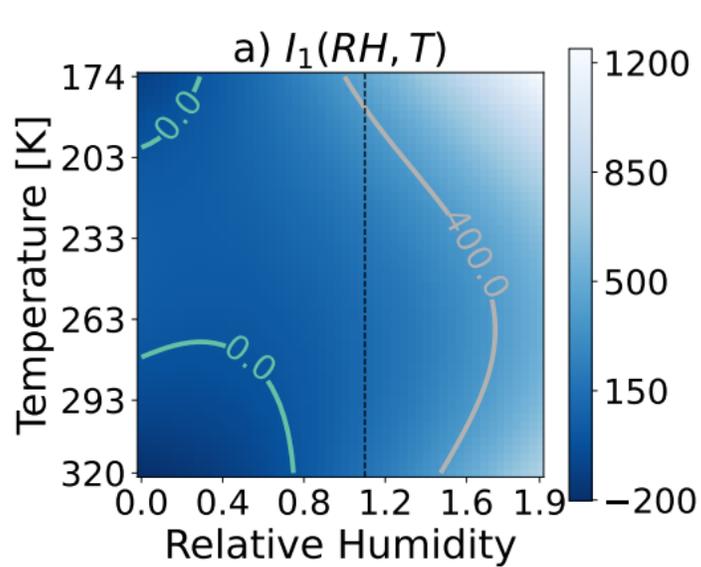


Figure 6.

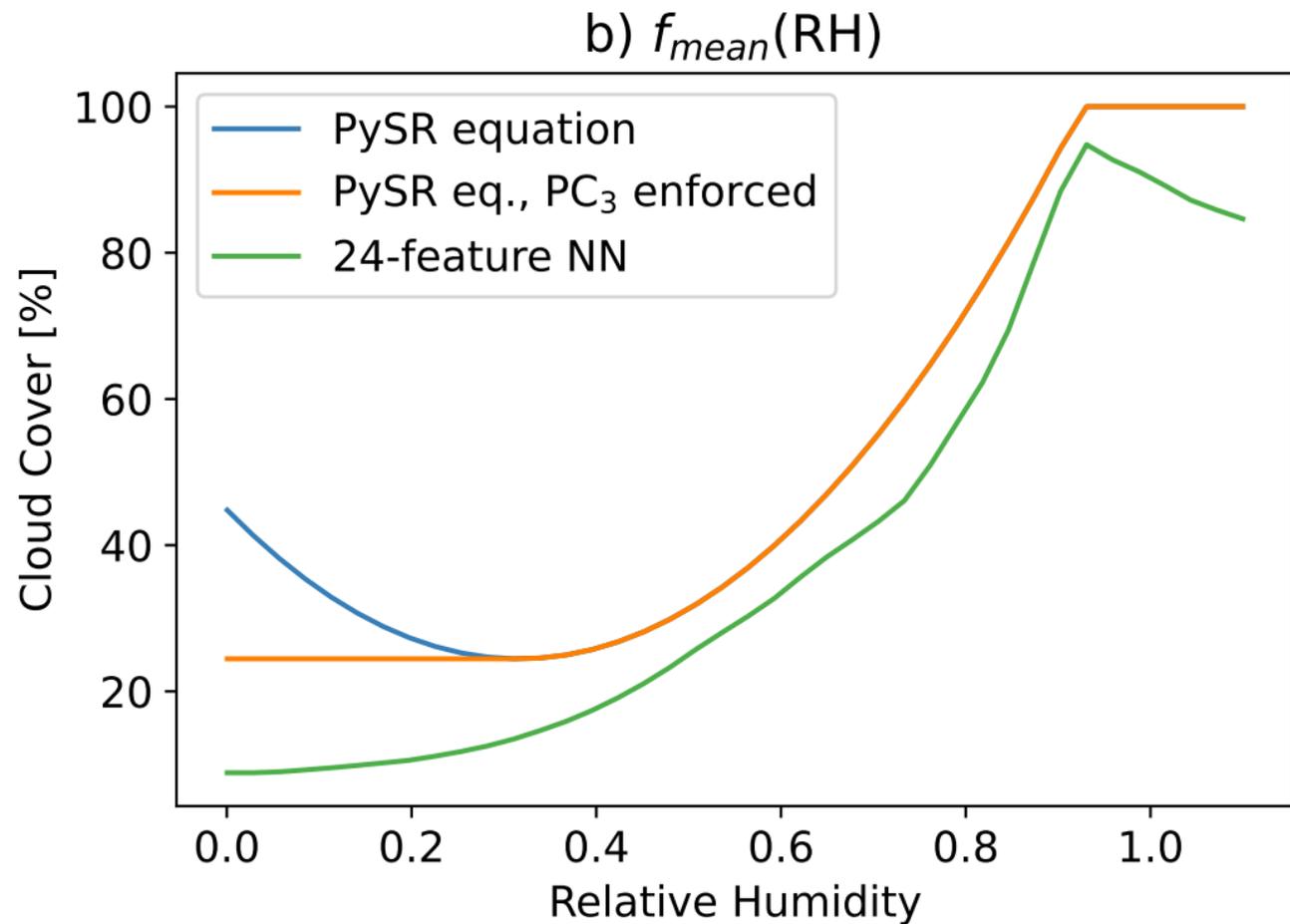
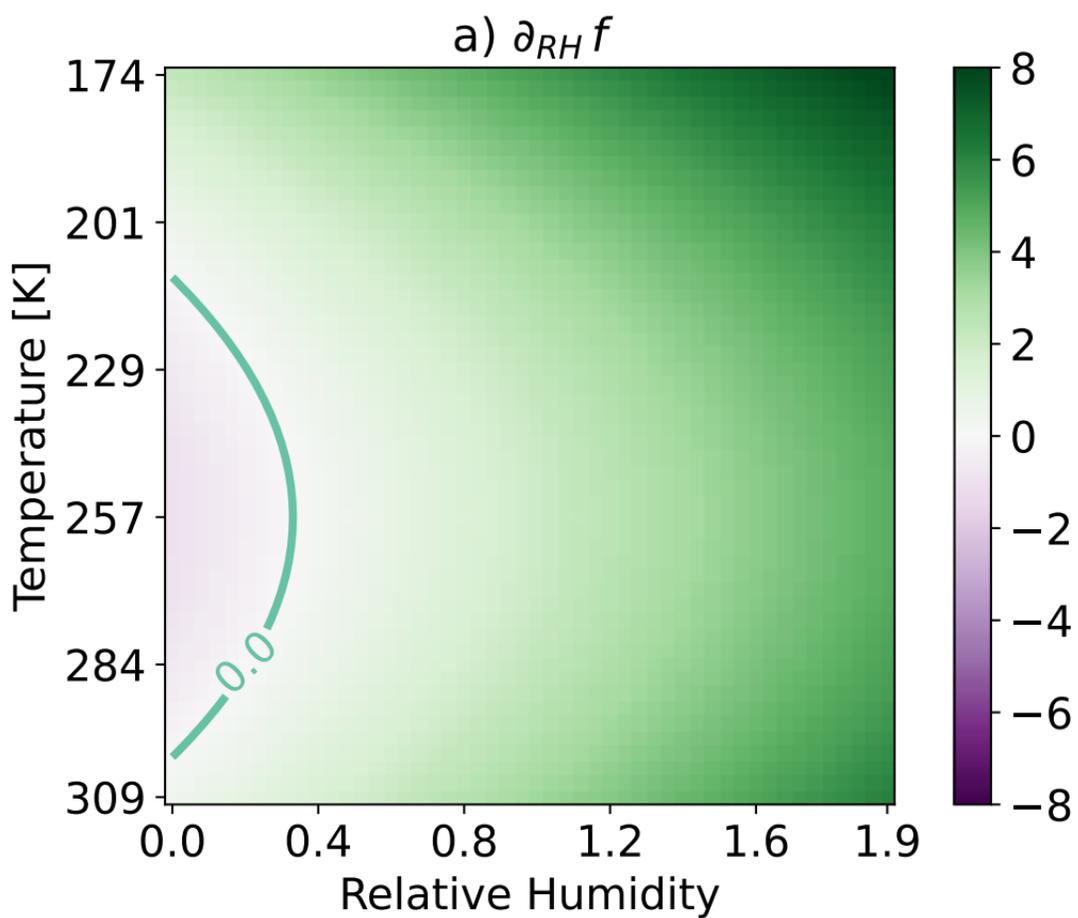


Figure 7.

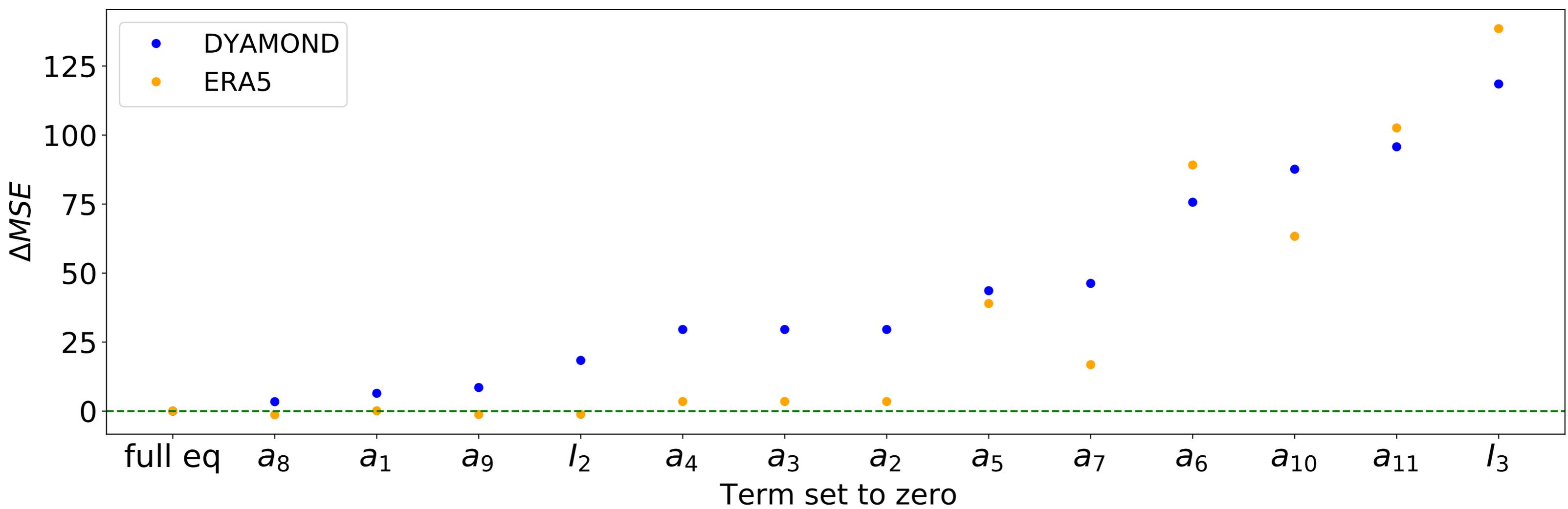
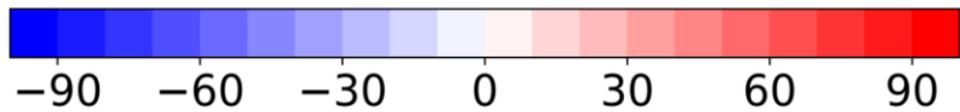
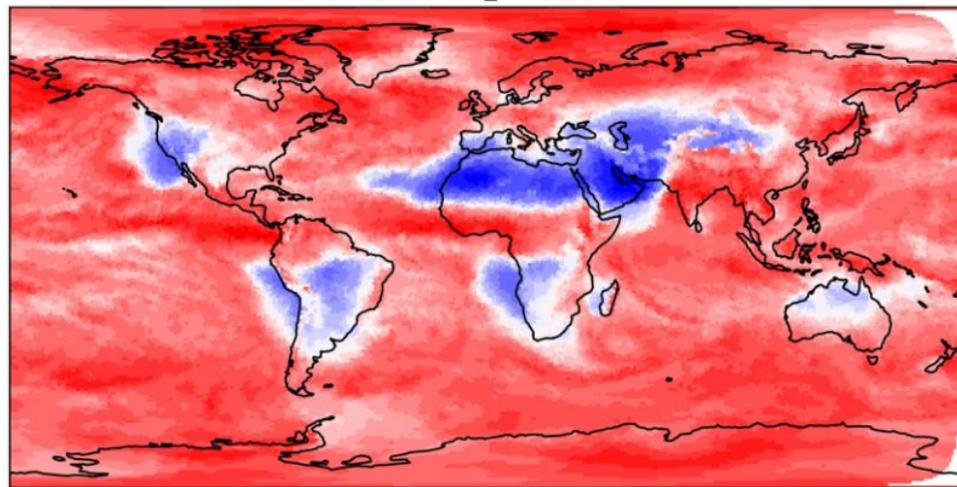
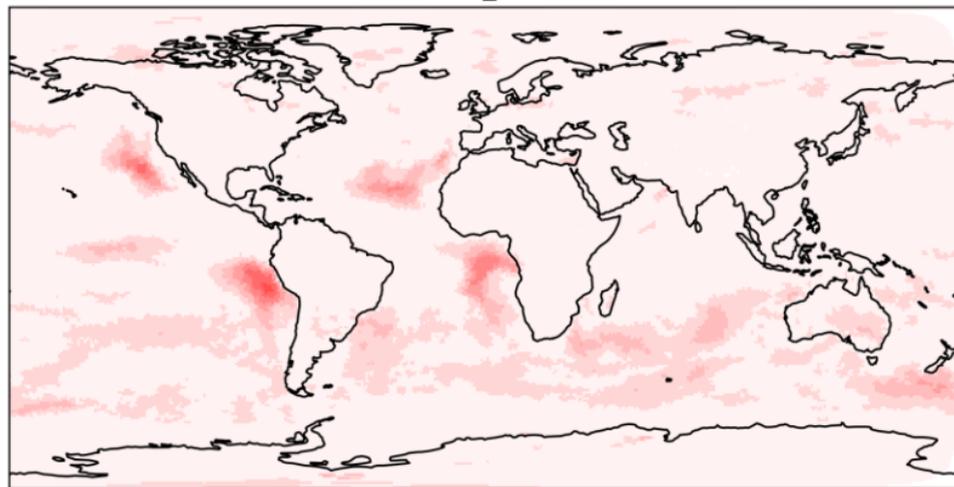


Figure 8.

l_1  l_2  l_3 