# Variant approach for identifying spurious relations that deep learning models learn

**Tobias Tesch** [1,2,*], **Stefan Kollet** [1,2] **and Jochen Garcke** [3,4]

[1] *Institute of Bio- and Geosciences, Agrosphere (IBG-3), Forschungszentrum Jülich, Jülich, Germany*

[2] *Center for High-Performance Scientific Computing in Terrestrial Systems, Geoverbund ABC/J, Jülich, Germany*

[3] *Fraunhofer Center for Machine Learning and Fraunhofer SCAI, Sankt Augustin, Germany*

[4] *Institut für Numerische Simulation, Universität Bonn, Bonn, Germany*

Correspondence*:
Tobias Tesch, Forschungszentrum Jülich IBG-3, Wilhelm-Johnen-Straße, 52428 Jülich, Germany
t.tesch@fz-juelich.de

## ABSTRACT

A deep learning (DL) model learns a function relating a set of input variables with a set of target variables. While the representation of this function in form of the DL model often lacks interpretability, several interpretation methods exist that provide descriptions of the function (e.g. measures of feature importance). On the one hand, these descriptions may build trust in the model or reveal its limitations. On the other hand, they may lead to new scientific understanding. In any case, a description is only useful if one is able to identify if parts of it reflect spurious instead of causal relations (e.g. random associations in the training data instead of associations due to a physical process). However, this can be challenging even for experts because, in scientific tasks, causal relations between input and target variables are often unknown or extremely complex. Commonly, this challenge is addressed by training separate instances of the considered model on random samples of the training set and identifying differences between the obtained descriptions. Here, we demonstrate that this may not be sufficient and propose to additionally consider more general modifications of the prediction task. We refer to the proposed approach as variant approach and demonstrate its usefulness and its superiority over pure sampling approaches with two illustrative prediction tasks from hydrometeorology. While being conceptually simple, to our knowledge the approach has not been formalized and systematically evaluated before.

Keywords: Interpretable deep learning, statistical model, machine learning, spurious correlation, causality, hydrometeorology, geoscience

## 1 INTRODUCTION

A deep learning (DL) model learns a function relating a set of input variables with a set of target variables. While DL models excel in terms of predictive performance, the representation of the learned function in form of the DL model (e.g. in form of a neural network) often lacks interpretability. To address this lack

of interpretability, several interpretation methods have been developed (see e.g. (Zhang and Zhu, 2018; Montavon et al., 2018; Molnar, 2019; Gilpin et al., 2018); (Samek et al., 2021)) providing descriptions of the learned function (e.g. measures of feature importance, FI). On the one hand, such descriptions can build trust in a model (Ribeiro et al., 2016) or reveal a model's limitations. Lapuschkin et al. (2019), for example, analyzed FI scores and found that their image classifier relied on a copyright tag on horse images. Similarly, Schramowski et al. (2020) analyzed FI scores and found (and corrected) that their DL model classified sugar beet leaves as healthy or diseased while incorrectly focusing on areas outside of the leaves.

On the other hand, descriptions of the learned function can lead to new scientific understanding. Ham et al. (2019), for example, analyzed FI scores and identified a previously unreported precursor of the Central-Pacific El Niño type; Gagne II et al. (2019) analyzed FI scores to gain a better understanding of the relations between environmental features and severe hail; McGovern et al. (2019) analyzed FI scores to gain a better understanding of the formation of tornadoes; and Toms et al. (2020) analyzed FI scores and identified regions related to the El Niño-Southern Oscillation (ENSO) and regions providing predictive capabilities for land surface temperatures at seasonal scales. Roscher et al. (2020) provide a general review of explainable machine learning for scientific insights in the natural sciences.

Whether descriptions of the function that a DL model learns are computed to build trust in the model, study the model's limitations, or gain new scientific understanding, it is important to identify if parts of a description reflect spurious instead of causal relations (e.g. random associations in the training data instead of associations due to a physical process). Examples for spurious relations are the above-mentioned copyright tag on horse images and the area outside of the classified sugar beet leaves. However, especially in prediction tasks involving physical, biological or chemical systems with several non-linearly interacting components, identifying spurious relations is challenging even for experts. Note that this does not only apply to the identification of spurious relations in descriptions of functions that DL models learn, but in general to the identification of spurious relations in descriptions of functions that any statistical model learns.

Commonly, this challenge is addressed by training separate instances of the considered model on random samples of the training set and aggregating or comparing the obtained descriptions. De Bin et al. (2015), for instance, compared subsampling and bootstrapping for the identification of relevant input variables in multivariable regression tasks. They applied a feature selection strategy repeatedly to samples of the original training set obtained by subsampling or bootstrapping, respectively, and identified relevant features by analyzing feature selection frequencies. As another example, Gagne II et al. (2019) trained 30 instances of different statistical models on sampled training and test sets to take into account that the models' skills and the relations between input and target variables that the models learn might depend on the specific training and test set composition. Here, we propose to not only consider sampling, but also more general modifications of the original prediction task. We refer to this more general approach as variant approach. In the approach, separate instances of the considered statistical model (referred to as variant models) are trained on modified prediction tasks (referred to as variant tasks) for which it is assumed that causal relations between input and target variables either remain stable or vary in specific ways. Subsequently, the descriptions of the functions that original and variant models learn are compared and it is evaluated whether they reflect the assumed stability or specific variation, respectively, of causal relations. If this is not the case for some parts of the descriptions, these parts likely reflect spurious relations. The approach constitutes a generalization of sampling approaches in that sampling is one of many ways for modifying the original prediction task in order to obtain a variant task.

68   A similar concept to ours has, to the best of our knowledge, only been pursued systematically in a strict
69   causality framework (for details on this framework see e.g. (Pearl, 2009) or for a more methodological
70   focus (Guo et al., 2020)). Peters et al. (2016), for example, consider modifications of an original prediction
71   task for which they require the conditional distribution $p(y|\vec{x}_S)$ of the target variable $y$ given the complete
72   set $\vec{x}_S$ of variables that directly cause $y$ to remain stable. Exploiting this requirement, they aim to identify
73   the subset $S$ of (direct) causal predictors within all observed features. While this approach is conceptually
74   related to the proposed variant approach, the latter does not require the strict causality framework but is
75   applicable to any machine learning prediction task. Note that in our work the terms causal and spurious do
76   not refer to an underlying causal graph or other concepts from the strict causality framework but should be
77   interpreted with common sense: a pixel in an image, for instance, is causally related to the label "dog" if
78   and only if it belongs to a dog in the image, and the value of a meteorological variable at a specific location
79   and time is causally related to the value of a meteorological variable at another location and time if and
80   only if one value influences the other via some physical process.

81   Other approaches in machine learning that consider modifications of an original prediction task
82   predominantly aim to improve the predictive performance of a statistical model rather than to analyze the
83   relations between input and target variables. Transfer learning (Pan and Yang, 2010), for instance, aims to
84   extract knowledge from one or more source tasks to apply it to a target task, e.g. training a neural network
85   first on a similar task before fine-tuning the weights on the target task. Adversarial training, as another
86   example, optimizes the loss over a set of perturbations of the input (Goodfellow et al., 2015; Sinha et al.,
87   2018) to become less susceptible to adversarial attacks (Szegedy et al., 2014), imperceptible changes to
88   the input that can change the model's prediction. Traditional importance weighting (Shimodaira, 2000) or
89   more recent methods (Lakkaraju et al., 2020), as further examples, shift the input distribution in order to
90   perform better on a known or unknown test distribution.

91   In this work, we demonstrate the proposed variant approach with two illustrative prediction tasks from
92   hydrometeorology. First, we predict the occurrence of rain at a target location, given geopotential fields
93   at different pressure levels in a surrounding region. Second, we predict the water level at a location in a
94   river, given the water level upstream and downstream 48 hours earlier. As statistical models, we consider
95   linear models and neural networks. After training a model on one of these tasks, we apply an interpretation
96   method to obtain a description of the learned function. This description indicates the average importance of
97   the different input locations for the predictions of the model. To identify if this importance reflects spurious
98   instead of causal relations between input and target variables, we apply the proposed variant approach.

99   The article is structured as follows: in Section Materials and Methods, we formalize the variant approach
100  and define the two prediction tasks and variants thereof that illustrate the approach. Further, we introduce
101  the statistical models and interpretation methods used in this work. Subsequently, we present and discuss
102  the results obtained when training the statistical models on the considered prediction tasks and applying
103  the variant approach. In Section Conclusions, we summarize our main findings and discuss perspectives for
104  future research and applications of the variant approach.

## 2   MATERIALS AND METHODS

### 2.1   Variant approach

106  During the training phase, a statistical model learns a function $f : \mathbb{R}^n \to \mathbb{R}^k$ relating an input space
107  $X \subseteq \mathbb{R}^n$ with a target space $Y \subseteq \mathbb{R}^k$ given a training set $T = \{(\vec{x_i}, \vec{y_i})\}_{i=1}^N$ with $\vec{x_i} \in X$, $\vec{y_i} \in Y$. As
108  the representation of $f$ in form of the statistical model (e.g. in form of a neural network) often lacks

109 interpretability, several interpretation methods have been developed (see e.g. (Zhang and Zhu, 2018;
110 Montavon et al., 2018; Molnar, 2019; Gilpin et al., 2018);(Samek et al., 2021)). Most of these methods
111 yield vector-valued descriptions $\vec{d} \in \mathbb{R}^d$ of $f$ (e.g. measures of feature importance). These descriptions can
112 be global or local, in the latter case not only depending on $f$ but on a subset $X_d \subset X$ as well. An example
113 of a global description are the weights of a linear regression model. An example of a local description $\vec{d}(\vec{x})$
114 is the gradient of a neural network evaluated at a location $\vec{x} \in X$.

115     A description $\vec{d}$ reflects the relations between input and target variables that the statistical model learned.
116 Whether the user aims to use $\vec{d}$ to build trust in the model, reveal the model's limitations, or gain new
117 scientific understanding, it is important to identify if parts of the vector $\vec{d}$ reflect spurious instead of causal
118 relations. In many cases, this is challenging even for experts. Therefore, we propose a variant approach.
119 The approach consists of three steps. First, the original prediction task is modified in such a way that
120 causal relations reflected in specific parts of $\vec{d}$ are assumed to either remain stable or vary in a specific
121 way. We refer to the modified prediction task as variant task. Second, a separate instance of the considered
122 statistical model (referred to as variant model) is trained on the variant task and a corresponding description
123 $\vec{d^v}$ (referred to as variant description) of the function $f^v$ that the variant model learns is computed. Third,
124 original and variant descriptions are compared and it is evaluated whether the previously specified parts of
125 original and variant descriptions reflect the assumed stability or specific variation, respectively, of causal
126 relations. If this is not the case, the respective parts of the vector $\vec{d}$ or of the vector $\vec{d^v}$ reflect spurious
127 relations.

128     Formalizing the approach, we define a variant task by an input space $X^v \subseteq \mathbb{R}^{n^v}$, a target space $Y^v \subseteq \mathbb{R}^{k^v}$,
129 a training set $T^v = \{(\vec{x_i^v}, \vec{y_i^v})\}_{i=1}^{N^v}$ with $\vec{x_i^v} \in X^v$, $\vec{y_i^v} \in Y^v$, an interpretation method (in most cases the
130 same as for the original task) that provides a description $\vec{d^v} \in \mathbb{R}^{d^v}$ of the learned function $f^v : \mathbb{R}^{n^v} \to \mathbb{R}^{k^v}$,
131 two sets of $m$ boolean vectors $\vec{I_j} \in \{0,1\}^d$ and $\vec{I_j^v} \in \{0,1\}^{d^v}$, $j = 1, \dots, m$, and $m$ corresponding smooth
132 (not necessarily symmetric) distance functions $dist_j : \mathbb{R}^d \times \mathbb{R}^{d^v} \to \mathbb{R}^{\geq 0}$, $j = 1, \dots, m$. We denote by
133 $\vec{d}(\vec{I_j})$ (and analogously by $\vec{d^v}(\vec{I_j^v})$) the restriction of $\vec{d}$ to the dimensions specified by the boolean vector $\vec{I_j}$
134 and refer to $\vec{d}(\vec{I_j})$ as a *part* of $\vec{d}$. The distance function $dist_j$ incorporates the user's assumption about how
135 the part $\vec{d}(\vec{I_j})$ of $\vec{d}$ changes for the variant task if it reflects causal relations, and quantifies the deviation of
136 this stability or specific variation, respectively. In other words, $dist_j$ computes a value $dist_j(\vec{d}, \vec{d^v})$ which is
137 0 if $\vec{d}(\vec{I_j})$ and $\vec{d^v}(\vec{I_j^v})$ exhibit the assumed stability or systematic variation, respectively, of causal relations.
138 In turn, the more they deviate from this assumed stability or specific variation, respectively, the larger the
139 value $dist_j(\vec{d}, \vec{d^v})$ should be.

140     Let us consider some examples of variant tasks. As already mentioned in the introduction, one way to
141 modify the original prediction task in order to obtain a variant task is to consider a sampled training set,
142 e.g. obtained by randomly sampling the original training set in the context of subsampling or bootstrapping
143 (De Bin et al., 2015). In this case, we assume that all causal relations remain stable. Hence, we may choose
144 to evaluate the dimensionwise distance between an original description $\vec{d} \in \mathbb{R}^d$ and the corresponding
145 variant description $\vec{d^v} \in \mathbb{R}^d$ of the function $f^v$ that a separate instance of the original model learns when
146 trained on the sampled training set. Using the above formalism, this corresponds to defining the boolean
147 vectors $(\vec{I_j})_i = (\vec{I_j^v})_i = \delta_{ji} \in \mathbb{R}^d$ (vectors with 0 components in all dimensions except from dimension $j$
148 where the component is 1) and the distance functions $dist_j(\vec{d}, \vec{d^v}) = |\vec{d}_j - \vec{d^v}_j|$ for $j = 1, \dots, m = d$. Now,
149 $dist_j(\vec{d}, \vec{d^v}) \gg 0$ for some $j \in \{1, \dots, d\}$ indicates that the part $\vec{d}(\vec{I_j}) = \vec{d}_j$ of the original description, or
150 the part $\vec{d^v}(\vec{I_j^v}) = \vec{d^v}_j$ of the variant description, reflects spurious relations. Note that we can repeat the
151 sampling procedure several times, leading to multiple variant tasks of the same type.

152    A second example for the definition of a variant task is to consider a modification of the input space.
153 Later, for instance, we consider the task to predict a rain event at a target location given input variables
154 in the $60 \times 60$ pixels neighborhood (see Fig. 1A). As a variant task, we consider the input variables in
155 the $80 \times 80$ pixels neighborhood instead. As original description $\vec{d} \in \mathbb{R}^{60 \times 60}$, we consider a measure
156 of the average importance of each pixel in the $60 \times 60$ pixels neighborhood for the predictions of the
157 original model, and as variant description $\vec{d^v} \in \mathbb{R}^{80 \times 80}$, we analogously measure the average importance
158 of each pixel in the $80 \times 80$ pixels neighborhood for the predictions of the variant model. In this case, we
159 assume that causal relations between pixels in the $60 \times 60$ pixels neighborhood and rain events at the target
160 location remain stable when enlarging the considered neighborhood by 10 pixels on each side. Hence, we
161 choose to evaluate the dimensionwise distance between the original description $\vec{d}$ and the central $60 \times 60$
162 pixels of the variant description $\vec{d^v}$. Using the above formalism, this corresponds to defining the boolean
163 matrices $(\vec{I}_{j_1 j_2})_{i_1 i_2} = \delta_{j_1 j_2, i_1 i_2} \in \mathbb{R}^{60 \times 60}$ (matrices with 0 components in all dimensions except from
164 dimension $j_1 j_2$ where the component is 1), the boolean matrices $(\vec{I^v}_{j_1 j_2})_{i_1 i_2} = \delta_{j_1 + 10 j_2 + 10, i_1 i_2} \in \mathbb{R}^{80 \times 80}$
165 (10 corresponds to the offset between the neighborhoods for original and variant task, i.e. input index
166 $(j_1 + 10, j_2 + 10)$ in the variant task corresponds to the same location as input index $(j_1, j_2)$ in the original
167 task) and the distance functions $dist_{j_1 j_2}(\vec{d}, \vec{d^v}) = |\vec{d}_{j_1 j_2} - \vec{d^v}_{j_1 + 10, j_2 + 10}|$ for $j_1, j_2 = 1, \ldots, 60$. Now,
168 $dist_{j_1 j_2}(\vec{d}, \vec{d^v}) \gg 0$ for some $j_1, j_2 \in \{1, \ldots, 60\}^2$ indicates that the part $\vec{d}(\vec{I}_{j_1 j_2}) = \vec{d}_{j_1 j_2}$ of the original
169 description, or the part $\vec{d^v}(\vec{I^v}_{j_1 j_2}) = \vec{d^v}_{j_1 + 10, j_2 + 10}$ of the variant description, reflects spurious relations.
170 Note that for some statistical models, this type of variant task might require slight changes to the model
171 architecture.

172    A third example for the definition of a variant task is to consider a modification of the target variable.
173 Later, for instance, we predict the water level at a location in a river given the water level in some specified
174 segment of the river (see Fig. 1B). As a variant task, we consider the same segment of the river but shift the
175 target location by $\tau$ pixels along the river (see Fig. 2B). As original and variant descriptions $\vec{d}, \vec{d^v} \in \mathbb{R}^d$, we
176 consider a measure of the average importance of each pixel in the specified river segment for the predictions
177 of the original model and the variant model, respectively. In this case, we assume that causal relations are
178 shifted along the river by the same distance as the target location is (i.e. by $\tau$ pixels). Hence, we choose
179 to compute the dimensionwise distance between the original description $\vec{d}$ and the variant description $\vec{d^v}$
180 shifted by $\tau$ dimensions (i.e. we consider the distance $|\vec{d}_j - \vec{d^v}_{j+\tau}|$ for all $j$ for that $j + \tau \in \{1, \ldots, d\}$).
181 Using the above formalism, this corresponds to defining the boolean vectors $(\vec{I}_j)_i = (\vec{I}_j^v)_{i+\tau} = \delta_{ji}$ and
182 the distance functions $dist_j(\vec{d}, \vec{d^v}) = |\vec{d}_j - \vec{d^v}_{j+\tau}|$ for all $j = 1, \ldots, d$ for that $j + \tau \in \{1, \ldots, d\}$. Now,
183 $dist_j(\vec{d}, \vec{d^v}) \gg 0$ indicates that the part $\vec{d}(\vec{I}_j) = \vec{d}_j$ of the original description, or the part $\vec{d^v}(\vec{I}_j^v) = \vec{d^v}_{j+\tau}$
184 of the variant description, reflects spurious relations.

185    In this example, it might be more realistic to assume that causal relations are not shifted along the river by
186 exactly $\tau$ pixels, but that the shift distance depends on the flow velocity and potentially further influences.
187 The proposed formalism allows to take this into account by varying the definition of $\vec{I}_j$, $\vec{I}_j^v$ and $dist_j$.
188 Suppose, for instance, that the flow velocity around the original target location is twice as high as around
189 the shifted target location. In this case, we might assume that the sum of importance of the *two* pixels
190 upstream of the original target location should be identical to the importance of the *single* pixel upstream
191 of the shifted target location. Hence, we might decide to consider $(\vec{I}_j)_i = \delta_{ji} + \delta_{j-1,i}$, $(\vec{I}_j^v)_{i+\tau} = \delta_{ji}$ (as
192 above), and $dist_j(\vec{d}, \vec{d^v}) = |(\vec{d}_j + \vec{d}_{j-1}) - \vec{d^v}_{j+\tau}|$, where the index $j$ corresponds to the original target
193 location. In this case, $dist_j(\vec{d}, \vec{d^v}) \gg 0$ indicates that the part $\vec{d}(\vec{I}_j)$ (corresponding to $\vec{d}_j$ *and* $\vec{d}_{j-1}$) of the
194 original description, or the part $\vec{d^v}(\vec{I}_j^v) = \vec{d^v}_{j+\tau}$ of the variant description, reflects spurious relations.

195  In general, however, it is difficult to take variations of flow velocity and further influences into account
196  when defining $\vec{I}_j$, $\vec{I}_j^v$ and $dist_j$. This is for example due to unavailable data on flow velocity and nonlinear
197  behavior (e.g. that the sum of importance of the *two* pixels upstream of the original target location should
198  be identical to the importance of the *single* pixel upstream of the shifted target location if the flow velocity
199  in the respective river segment is twice as high, likely represents a too strong assumption on linearity). We
200  will come back to this in the discussion of the results.

201  Let us return to the formal definition of the variant approach. The first step was to define a variant task.
202  The second step consists of training a separate instance of the original model (a variant model) on this
203  task and computing a variant description. The third step of the approach consists of comparing original
204  and variant description and evaluating $dist_j(\vec{d}, \vec{d^v}) \gg 0$ for all $j = 1, \ldots, m$. If $dist_j(\vec{d}, \vec{d^v}) \gg 0$ for some
205  $j \in \{1, \ldots, m\}$, the user infers that $\vec{d}(\vec{I}_j)$ or $\vec{d^v}(\vec{I}_j^v)$ reflects spurious relations. Note that the converse is
206  not possible, i.e. if $dist_j(\vec{d}, \vec{d^v}) \approx 0$, the user cannot infer that $\vec{d}(\vec{I}_j)$ reflects causal relations (as it might be
207  that both $\vec{d}(\vec{I}_j)$ and $\vec{d^v}(\vec{I}_j^v)$ reflect spurious relations). Note further that the specification of the condition
208  $dist_j(\vec{d}, \vec{d^v}) \gg 0$ should in general take into account the specific original and variant task, the choice of
209  the distance function $dist_j$, and the certainty of the assumed stability or systematic variation, respectively,
210  of causal relations. Moreover, in case the user does not need a binary identification of parts of $\vec{d}$ that reflect
211  spurious relations, it might be better not to consider the binary condition $dist_j(\vec{d}, \vec{d^v}) \gg 0$, but to consider
212  raw values $dist_j(\vec{d}, \vec{d^v})$, where higher distances indicate a higher probability that $\vec{d}(\vec{I}_j)$ or $\vec{d^v}(\vec{I}_j^v)$ reflects
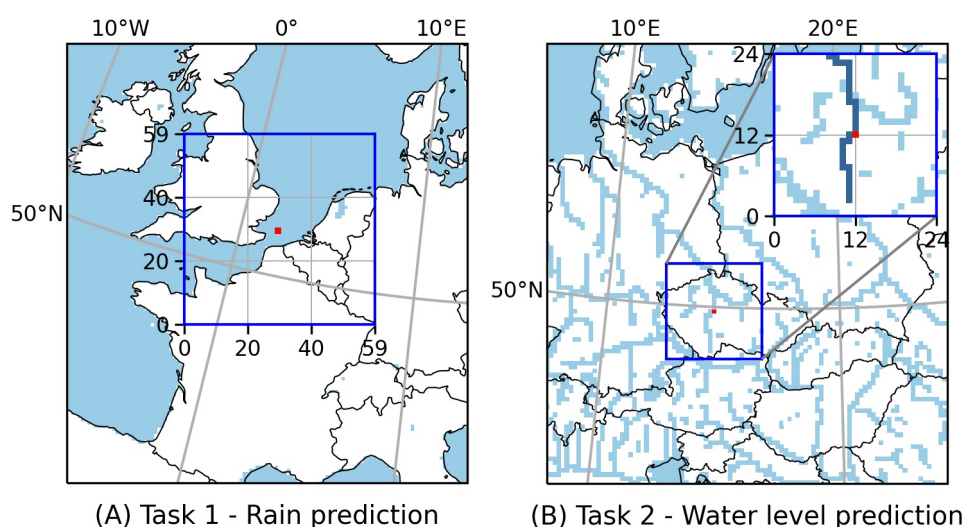213  spurious relations.

214  For all variant tasks defined in this work, the expression $dist_j(\vec{d}, \vec{d^v})$ corresponds to the relative distance
215  between a single component $\vec{d}_{j_1}$ of an original description and a single component $\vec{d^v}_{j_2}$ of a corresponding
216  variant description, i.e. it takes the form

$$dist_j(\vec{d}, \vec{d^v}) = \frac{|\vec{d}_{j_1} - \vec{d^v}_{j_2}|}{|\vec{d}_{j_1}| + |\vec{d^v}_{j_2}| + \varepsilon}, \tag{1}$$

217  with some regularization parameter $\varepsilon \geq 0$. By considering relative distances rather than absolute distances,
218  we define, for instance, that $\vec{d}_{j_1} = 100$, $\vec{d^v}_{j_2} = 101$ agree better than $\vec{d}_{j_1} = 1$, $\vec{d^v}_{j_2} = 2$, or, in other words,
219  in the latter case it is more likely that the value $\vec{d}_{j_1}$ or the value $\vec{d^v}_{j_2}$ reflects spurious relations. Further, an
220  advantage of considering relative distances is that all distances lie between zero and one (when neglecting
221  $\varepsilon$) which allows to apply a threshold $t \in (0, 1)$ to specify the condition $dist_j(\vec{d}, \vec{d^v}) \gg 0$ and to mark all
222  parts $\vec{d}(I_j)$ of the original description as spurious for which $dist_j(\vec{d}, \vec{d^v}) > t$. In this study, we use $t = 0.5$
223  as threshold and $\varepsilon = 1e - 3$ as regularization parameter. Choosing a smaller threshold, more values are
224  marked as spurious (with all values marked as spurious for $t = 0$), and choosing a larger threshold, fewer
225  values are marked as spurious (with no values marked as spurious for $t = 1$) by definition. For the examples
226  considered below, $t = 0.5$ seems to be a good choice.

## 227  2.2  Illustrative tasks

228  In this section, we define two prediction tasks and corresponding variant tasks that illustrate the proposed
229  variant approach. We chose simplified tasks and global descriptions of the learned functions to be able
230  to decide whether parts of the descriptions that the variant approach marks as spurious do indeed reflect
231  spurious relations. The data underlying both tasks is 3-hourly data at $412 \times 424$ pixels over Europe. The
232  data was obtained from a long-term (January 1996 - August 2018), high-resolution ($\approx 12.5$ km) simulation
233  (Furusho-Percot et al., 2019) performed with the Terrestrial Systems Modeling Platform (TSMP), a fully

(A) Task 1 - Rain prediction      (B) Task 2 - Water level prediction

**Figure 1.** Set up of the two original prediction tasks. **(A)** Predict whether the precipitation averaged over the red $2 \times 2$ pixels target patch in the center of the $60 \times 60$ pixels input region exceeds 1 mm in the next 3 hours. **(B)** Predict the water level at the red pixel given the water level 48 hours earlier at the red pixel and the pixels upstream and downstream marked dark blue in the inset. Light blue indicates pixels with ponded water at the land surface during the entire simulation period (rivers, lakes, . . . ).

234    integrated groundwater-soil-vegetation-atmosphere modeling system (Shrestha et al., 2014; Gasper et al.,
235    2014). Note that the statistical models and interpretation methods applied in this work are described in
236    Section Statistical Models and Descriptions.
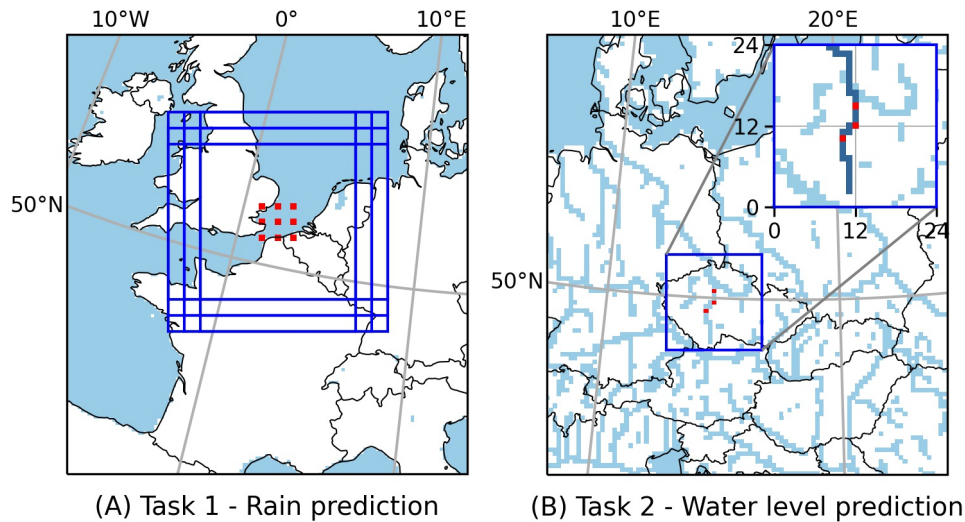
### 2.2.1    Task 1 – Rain prediction

238    In the first example, we predict the occurrence of rain at a $2 \times 2$ pixels target patch, given the geopotential
239    fields at 500, 850 and 1000 hPa in the $60 \times 60$ pixels neighborhood (see Fig. 1A). We model this as a
240    classification task and define that rain occurred, if the precipitation averaged over the target patch exceeds
241    1 mm in the following 3 hours. Previous works (Larraondo et al., 2019; Pan et al., 2019) have used CNNs to
242    predict precipitation given geopotential fields to improve the parameterization of precipitation in numerical
243    weather prediction models. Thus, apart from the simplifications of only one target location and a binary
244    target, this is a realistic prediction task.

245    As statistical models, we consider a logistic regression model and two convolutional neural networks
246    (CNNs) of different depth and complexity. As description of the function that the logistic regression
247    model learns, we consider the absolute values of the model weights averaged over the pressure level
248    axis. As descriptions of the functions that the CNNs learn, we consider saliency maps averaged over the
249    pressure level axis and over all training samples. These descriptions can be seen as measures of the average
250    importance of each pixel in the $60 \times 60$ pixels input region for the predictions of the models (for details see
251    the respective sections below).

252    To identify whether parts of the descriptions reflect spurious relations that the models learned, we compute
253    descriptions for variant models trained on three types of variant tasks. The first type (later referred to as
254    sampling type) considers the same task, but a modified training set obtained by randomly sampling 70 %
255    of the original training set without replacement. In this case, we assume that all causal relations remain

256 stable. Hence, we compute the pixelwise distance between original and variant descriptions. We repeat the
257 sampling procedure 10 times obtaining 10 variant tasks of this type. The second type of variant tasks (later
258 referred to as size type) considers the same task but the input variables in the $80 \times 80$ pixels neighborhood
259 of the target patch. In this case, we assume that causal relations between pixels in the $60 \times 60$ pixels
260 neighborhood and rain events at the target patch remain stable when enlarging the considered neighborhood
261 by 10 pixels on each side. Hence, we compute the pixelwise distance between the original descriptions and
262 the central $60 \times 60$ pixels of the variant descriptions. The third type of variant task (later referred to as
263 location type) considers the same task but for eight different target patches obtained by moving the original
264 target patch by five pixels to the left or right, and up or down. The input regions are shifted accordingly (see
265 Fig. 2A). In this case, we assume again that all causal relations remain stable. Hence, we again compute
266 the pixelwise distance between original and variant descriptions.

267   Note that to compute the variant descriptions for the functions that separate instances of the CNNs learn
268 when trained for different target locations, we average the saliency maps over all training samples from
269 the *original* task. This is because the distribution $p(\vec{x})$ of geopotential fields differs at different locations.
270 Thus, if we averaged the saliency maps for a variant CNN over all training samples from a variant task, the
271 obtained variant description would differ from the original description even if original and variant models
272 learned the exact same function relating geopotential fields and rain events.



(A) Task 1 - Rain prediction          (B) Task 2 - Water level prediction

**Figure 2.** Location variant tasks. **(A)** Original target patch (center) with its input region and eight additional
target patches and their (overlapping) input regions. **(B)** Original target location (center) and two additional
target locations closely upstream and downstream.

273   We obtained the geopotential fields and precipitation data from the aforementioned simulation. We
274 selected the geopotential fields in the considered input regions and created the binary rain event time series
275 for the corresponding target patches. Next, we split the time series using the first 56,000 time steps as
276 training candidates and the last 10,183 time steps as validation candidates. Finally, training and validation
277 sets were obtained by selecting all time steps followed by a rain event at the considered target patch and an
278 equal amount of randomly chosen additional time steps for non-rain events from the training and validation
279 candidates, respectively. This resulted in balanced training and validation sets of a total of approximately

280  10,000 time steps for each target patch. Handling strongly unbalanced data sets as it would be necessary
281  without such a selection of time steps is out of scope for this work.

### 2.2.2　Task 2 – Water level prediction

283  As a second example, we predict the water level at a location in a river, given the water level in a specific
284  segment of the river 48 hours earlier (see Fig. 1B).

285  As statistical models, we consider a linear regression model and a multilayer perceptron (MLP). As
286  description of the function that the linear regression model learns, we consider as in Task 1 the absolute
287  values of the model weights. For the MLP, we consider again the saliency maps averaged over all training
288  samples. Analogously to Task 1, these descriptions can be seen as measures of the average importance of
289  each pixel in the considered river segment for the predictions of the models (for details see the respective
290  sections below).

291  To identify whether parts of the descriptions reflect spurious relations that the models learned, we compute
292  descriptions for variant models trained on two types of variant tasks. The first type (later referred to as
293  sampling type) considers the same task, but a modified training set obtained by randomly sampling 70 %
294  of the original training set without replacement. In this case, we assume that all causal relations remain
295  stable. Hence, we compute the pixelwise distance between original and variant descriptions. We repeat the
296  sampling procedure 10 times obtaining 10 variant tasks of this type. The second type of variant tasks (later
297  referred to as location type) considers the same river segment as input, but target locations closely upstream
298  and downstream of the original target location (see Fig. 2B). In this case, we assume that causal relations
299  are shifted along the river by the same distance as the target location is. Hence, we compute the pixelwise
300  distance between the original description $\vec{d}$ and the variant description $\vec{d^v}$ shifted by $\tau$ pixels, where $\tau$ is
301  the number of pixels that the target location was shifted (i.e. we consider the distance $|\vec{d}_j - \vec{d^v}_{j+\tau}|$ for all $j$
302  for that $j + \tau \in \{1, \ldots, d\}$).

303  We obtained the water level data from the aforementioned simulation. In contrast to Task 1, this task is
304  not a classification but a regression task; discarding time steps to obtain a balanced data set is not necessary.
305  Hence, we use water level data for all 64,240 3-hourly time steps between January 1996 and December
306  2017. We randomly selected the years 1997, 2004, 2008 and 2015 as test data, covering the whole period
307  of time, and use the remaining years to train the models.

## 2.3　Statistical models and descriptions

309  In this section, we present the statistical models used in this study. Further, we describe saliency maps,
310  the interpretation method applied to obtain descriptions of the functions that the neural networks (MLP
311  and CNNs) learn. Note that for the considered examples, layerwise relevance propagation (LRP) and
312  Grad-CAM give very similar results to saliency maps. The section is ordered with respect to the complexity
313  of the described methods from simple to complex.

### 2.3.1　Linear Regression

315  Given training samples $(\vec{x}_i, y_i)_{i=1}^n$ with $\vec{x}_i \in \mathbb{R}^N$, $y_i \in \mathbb{R}$, a linear regression model learns a function
316  $f : \mathbb{R}^N \to \mathbb{R}$ of the form

$$f(\vec{x}) = \beta_0 + \vec{x}^T \cdot \vec{\beta}, \tag{2}$$

317   where $\vec{\beta} = (\beta_0, \vec{\vec{\beta}}) = (\beta_0, \beta_1, \ldots, \beta_N) \in \mathbb{R}^{N+1}$ are the weights of the model. Those weights are obtained
318   by minimizing the squared error on the training set

$$\sum_{i=1}^{n} (f(\vec{x}_i) - y_i)^2. \tag{3}$$

319   Optionally, a regularization term can be added to the objective. We calculate the minimizing weights $\vec{\beta}$
320   using the implementation of scikit-learn (Pedregosa et al., 2011). In our case, the inputs $\vec{x}_i$ are elements
321   of $\mathbb{R}^{30}$ representing the water level at the 30 pixels in the considered river segment (see Fig. 1B) and the
322   targets $y_i \in \mathbb{R}$ represent the water level at the target pixel 48 hours later.

323   As description of the function that a linear regression model learned, we consider the absolute values of
324   the weights $\vec{\vec{\beta}}$. This can be seen as a measure of the average importance of each pixel in the river segment
325   for the predictions of the model (Molnar, 2019).

### 2.3.2   Logistic Regression

327   Given the task to predict a binary target $y \in \{0, 1\}$ from an input $\vec{x} \in \mathbb{R}^N$, a logistic regression model
328   yields

$$P(y = 1 | \vec{x}, \vec{\beta}) = \frac{1}{1 + \exp(-(\beta_0 + \vec{x}^T \cdot \vec{\vec{\beta}}))}, \tag{4}$$

329   where $\vec{\beta} = (\beta_0, \vec{\vec{\beta}}) = (\beta_0, \beta_1, \ldots, \beta_N) \in \mathbb{R}^{N+1}$ are the weights of the model. These weights are obtained
330   by minimizing the function

$$-\prod_{i=1}^{n} P(y_i = 1 | \vec{x}_i, \vec{\beta})^{y_i} \cdot (1 - P(y_i = 1 | \vec{x}_i, \vec{\beta}))^{1-y_i} + \lambda R(\vec{\beta}) \tag{5}$$

331   with respect to $\vec{\beta}$. Here, $(\vec{x}_i, y_i)_{i=1}^n$ are training samples with $\vec{x}_i \in \mathbb{R}^N, y_i \in \{0, 1\}$, and $\lambda R(\vec{\beta})$ is a
332   regularization term. The product represents the probability with that – according to the logistic regression
333   model with weights $\vec{\beta}$ – the targets $y_i$ are observed given the input samples $\vec{x}_i$. Thus, minimizing the
334   negative product with respect to $\vec{\beta}$ corresponds to finding the $\vec{\beta}$ for that the highest probability is assigned
335   to observing the targets $y_i$ given the inputs $\vec{x}_i$ from the training set. We use scikit-learn (Pedregosa et al.,
336   2011) (solver 'liblinear') to approximate the minimizing weights $\vec{\beta}$. In our case, the inputs $\vec{x}_i$ are the
337   geopotential fields at 500, 850 and 1000 hPa flattened to vectors in $\mathbb{R}^{3 \cdot 60 \cdot 60}$ and the targets $y_i \in \{0, 1\}$
338   represent whether a rain event took place or not.

339   As description of the function that a logistic regression model learned, we consider the weights $\vec{\vec{\beta}}$. We
340   reshape the vector $\vec{\vec{\beta}}$ to the shape of the original input, $3 \times 60 \times 60$, take the absolute value and build an
341   average over the first (pressure level) axis. This can be seen as a measure of the average importance of each
342   pixel in the $60 \times 60$ pixels input region for the predictions of the model (Molnar, 2019).

### 2.3.3   Multilayer Perceptron

344   Multilayer Perceptrons (MLPs), also referred to as fully-connected neural networks, are feedforward
345   artificial neural networks. They are composed of one or more hidden layers and an output layer. Each
346   layer comprises several neurons. Each neuron in the first hidden layer builds a weighted sum of all
347   input variables, while each neuron in the subsequent layers builds a weighted sum of the outputs of the

348 neurons in the respective previous layer. In case of a neuron in a hidden layer, the sum is passed through a
349 nonlinear activation function and forms the input to the next layer. In case of a neuron in the output layer,
350 the sum is optionally passed through a nonlinear activation function and forms the output of the neural
351 network. The weights of the MLP are learned by minimizing a loss function on training samples $(\vec{x}_i, \vec{y}_i)_{i=1}^n$,
352 $\vec{x}_i \in \mathbb{R}^N, \vec{y}_i \in \mathbb{R}^K$, using backpropagation (LeCun et al., 2012).

353     In our case, the inputs to the MLP are elements $\vec{x}$ of $\mathbb{R}^{30}$ representing the water level at the 30 pixels
354 in the considered river segment (see Fig. 1b). The targets $y_i \in \mathbb{R}$ represent the water level at the target
355 pixel 48 hours later. Section Saliency maps describes how we obtained a description of the function that
356 the MLP learned. The network and training of the MLP were implemented using the deep learning library
357 Pytorch (Paszke et al., 2019). A detailed description of the used architecture and training procedure can be
358 found in the Supplementary Information.

### 2.3.4 Convolutional Neural Networks

360     Convolutional Neural Networks (CNNs) are frequently employed DL models designed to process stacks
361 of multiple arrays containing spatially structured data. This can, for example, be a stack of 2-dimensional
362 arrays for an RGB image ($\vec{x}_i \in \mathbb{R}^{3 \times \text{height} \times \text{width}}$) or, as in our case, a stack of 2-dimensional geopotential
363 fields at different pressure levels in the atmosphere ($\vec{x}_i \in \mathbb{R}^{3 \times 60 \times 60}$). Typically, a CNN consists of three
364 types of layers: convolutional layers, pooling layers and fully-connected layers. In the following short
365 review of the typical CNN layers, we consider the case of one or multiple 2-dimensional input arrays. A
366 generalization of the concepts to N-dimensional input arrays is straightforward.

367     The input to a convolutional layer is a stack of $c_{\text{in}}$ 2-dimensional arrays and its output is a stack of $c_{\text{out}}$
368 2-dimensional arrays. The convolutional layer is characterized by $c_{\text{out}}$ kernels, which are 3-dimensional
369 tensors of shape $c_{\text{in}} \times k \times k$, where the kernel size $k$ is usually between 1 and 7. The output of the layer
370 are the $c_{\text{out}}$ 2-dimensional arrays obtained by convolving the input with each kernel along the last two
371 dimensions. Usually, a convolutional layer is directly followed by a nonlinear activation function which
372 is applied elementwise to the layer's output. In contrast to a fully-connected layer, a convolutional layer
373 preserves the spatial structure of the input: only neurons in a neighborhood defined by the kernel size
374 influence the output of a specific neuron.

375     As for convolutional layers, the input to a pooling layer is a stack of $c_{\text{in}}$ 2-dimensional arrays of shape
376 $n \times m$. Pooling layers reduce the dimensionality of the 2-dimensional arrays creating invariances to small
377 shifts and distortions. A typical form of pooling is max-pooling with a kernel size of two. This reduces the
378 resolution along both axes of each of the $c_{\text{in}}$ 2-dimensional arrays by a factor of two, picking always the
379 maximum value of a $2 \times 2$ patch of the original array. Thus, the output of this pooling layer is a stack of
380 $c_{\text{out}} = c_{\text{in}}$ 2-dimensional arrays of shape $\frac{n}{2} \times \frac{m}{2}$.

381     After several alternating convolutional and pooling layers which extract features of increasing complexity,
382 the resulting $c$ 2-dimensional arrays are flattened into a single vector and one or more fully-connected
383 layers, as described for the MLP, follow. The weights for the kernels in the convolutional layers and
384 the fully-connected layers are learned by minimizing a loss function on training samples $(\vec{x}_i, \vec{y}_i)_{i=1}^n$,
385 $\vec{x}_i \in \mathbb{R}^N, \vec{y}_i \in \mathbb{R}^K$, using backpropagation (LeCun et al., 2012). To prevent CNNs from overfitting,
386 dropout regularization (Srivastava et al., 2014) and batch normalization (Ioffe and Szegedy, 2015) are
387 commonly employed techniques.

388     In our case, the inputs $\vec{x}_i$ are the geopotential fields at 500, 850 and 1000 hPa, $\vec{x}_i \in \mathbb{R}^{3 \times 60 \times 60}$. The targets
389 $y_i \in \{0, 1\}$ represent whether a rain event took place or not. We consider two convolutional neural networks

390  of different depth and complexity. CNN1 is a shallow CNN with only two convolutional layers followed
391  by a single fully-connected layer. CNN2 is a commonly employed, much deeper CNN architecture called
392  resnet18 (He et al., 2016) for which the last fully-connected layer was adapted to have only two output
393  neurons to fit our binary prediction task. Section Saliency maps describes how we obtained descriptions of
394  the functions that the CNNs learned. The networks and training were implemented using the deep learning
395  library Pytorch (Paszke et al., 2019). A detailed description of the used CNN architectures and training
396  procedure can be found in the Supplementary Information.

### 2.3.5 Saliency maps

398  A common subgroup of interpretation methods providing descriptions of the functions that neural
399  networks (NNs) learn, are methods that assign an importance to each dimension of individual input samples
400  $\vec{x} \in \mathbb{R}^N$ (local feature importance scores), see e.g. (Samek et al., 2021). Among the most employed and
401  well-known methods for that purpose are saliency maps (Simonyan et al., 2013), layerwise relevance
402  propagation (LRP) (Bach et al., 2015) and Grad-CAM (Selvaraju et al., 2017). In the examples presented
403  in this work, all three methods yield similar results. Therefore and for the sake of brevity, we focus on
404  saliency maps (although e.g. Montavon et al. (2018) argue that saliency maps provide a bad measure of
405  feature importance because they indicate how the prediction of a model changes when the value of a feature
406  is changed, rather than indicating what makes the model make a prediction).

407  Note that in contrast to the weights of linear and logistic regression models, saliency maps are local
408  descriptions of the learned functions, i.e. the importance assigned to an input dimension (in our case an
409  input pixel) depends on the input sample $\vec{x}$. To get a global description of the learned function and a
410  measure of the average importance of each input pixel, we average the saliency maps over all training
411  samples.

412  In the rain prediction task, the NN defines an (almost everywhere) differentiable function $f$ that maps
413  geopotential fields $\vec{x} \in \mathbb{R}^{3 \times 60 \times 60}$ to probabilities $f(\vec{x}) = y \in (0, 1)$ that a rain event occurs. The partial
414  derivative

$$w_{cij}(\vec{x}) = \frac{\partial f}{\partial x_{cij}}(\vec{x}), \ \ c = 1, 2, 3, \ \ i, j = 1, \ldots, 60 \tag{6}$$

415  indicates how a small perturbation of the $c$-th geopotential field at pixel $(i, j)$ affects the prediction of the
416  NN. The saliency map

$$M_{ij}(\vec{x}) = \frac{1}{3} \sum_{c=1}^{3} |w_{cij}(\vec{x})|, \ \ i, j = 1, \ldots, 60 \tag{7}$$

417  considers the absolute value of the partial derivatives averaged over the pressure level axis to obtain for
418  each pixel in the $60 \times 60$ pixels input region a measure of its importance for the model's prediction for
419  sample $\vec{x}$.

420  In the water level prediction task, the neural network maps water levels $\vec{x} \in \mathbb{R}^{30}$ to a water level prediction
421  $f(\vec{x}) = y \in \mathbb{R}$. The saliency map

$$M_i(\vec{x}) = |w_i(\vec{x})| = \left| \frac{\partial f}{\partial x_i}(\vec{x}) \right|, \ \ i = 1, \ldots, 30 \tag{8}$$
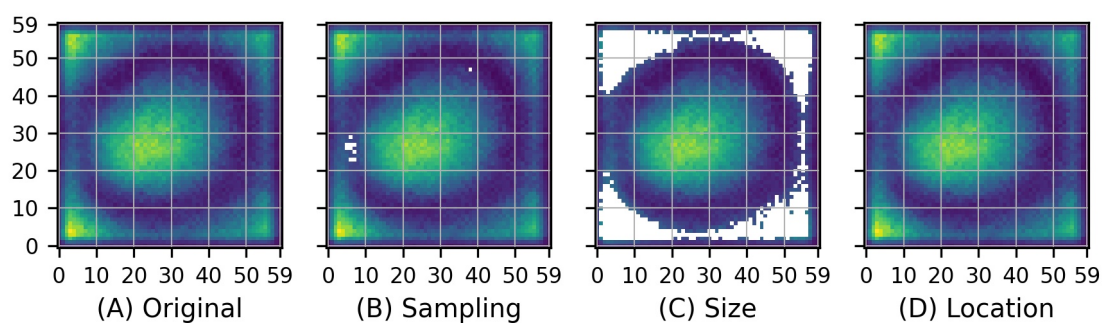
422  provides for each pixel in the considered river segment a measure of its importance for the model's
423  prediction for sample $\vec{x}$.

## 3 RESULTS AND DISCUSSION

### 3.1 Task 1 – Rain prediction

Figure 3A shows the description $\vec{d}$ of the function that CNN1 learned when it was trained on the original rain prediction task. Remember that the considered description is a measure of the average importance of each pixel in the $60 \times 60$ pixels input region for the predictions of the model. Our objective is to apply the variant approach to identify parts of the description that reflect spurious relations. To that purpose, we defined several variant tasks above. As a next step, we computed the corresponding variant descriptions, i.e. the descriptions of the functions that separate instances of CNN1 learned when trained on these variant tasks. For illustration, Figure 4 shows the original description (center, same as Fig. 3A) and the variant descriptions $\vec{d^{v_i}}, i = 1, \ldots, 8$, obtained for the eight location variant tasks (see Fig. 2A).

For each of these variant descriptions $\vec{d^{v_i}} \in \mathbb{R}^{60 \times 60}, i = 1, \ldots, 8$, we evaluated the pixelwise relative distance to the original description $\vec{d} \in \mathbb{R}^{60 \times 60}$ (see Equation 1), and masked all pixels of the original description $\vec{d}$ for which this distance exceeds the threshold of $t = 0.5$ for any $\vec{d^{v_i}}$. The resulting masked version of $\vec{d}$ is shown in Fig. 3D. Note that in this case, there is no pixel for which the relative distance between original description and any of the variant descriptions exceeds 0.5, hence Fig. 3D is identical to Fig. 3A. Analogously to Fig. 3D, Fig. 3B shows the masked version of $\vec{d}$ obtained when masking all pixels for which the pixelwise relative distance between $\vec{d}$ and one of the variant descriptions $\vec{d^{v_i}}$ obtained for the sampling variant tasks exceeds 0.5. We observe that some pixels in the west of the inner area of importance are masked, indicating that the inner area of importance might actually extend further to the west. Figure 3C shows the masked version of $\vec{d}$ obtained when masking all pixels for which the pixelwise relative distance between $\vec{d}$ and the central $60 \times 60$ pixels of the variant description $\vec{d^{v_i}}$ obtained for the size variant task exceeds 0.5. Notably, all the boundary pixels with high values in Fig. 3a are masked, indicating that these values likely reflect spurious relations.



**Figure 3.** **(A)** Description $\vec{d}$ of the function that CNN1 learned when it was trained on the original rain prediction task (see Fig. 1a). The description is a measure of the average importance of each pixel in the $60 \times 60$ pixels input region for the predictions of the model. Yellow color indicates high and blue color low importance. **(B-D)** As (A), but pixels for which the relative distance between the original description $\vec{d}$ and one of the variant descriptions $\vec{d^{v_i}}$ obtained for the sampling, size and location variant tasks, respectively, exceeds the threshold of $t = 0.5$, are masked.

Figure 5 shows the same as Fig. 3 but for CNN2. Only few pixels are masked for the sampling and location variant tasks. However, the mask obtained for the size variant task indicates that the checkerboard pattern in the original description $\vec{d}$, which is shown in Fig. 5A, likely reflects spurious relations. Note that
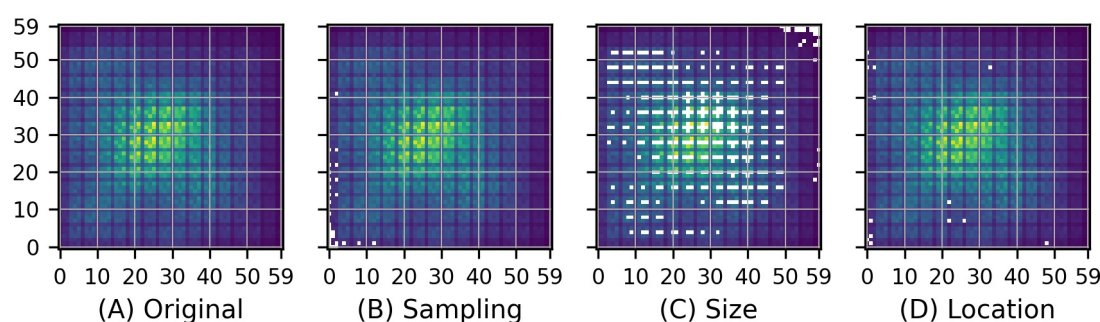
**Figure 4.** Descriptions obtained when training separate instances of CNN1 for the nine different locations depicted in Fig. 2A. Each description is a measure of the average importance of each pixel in the $60 \times 60$ pixels input region for the predictions of the respective instance of CNN1. Yellow color indicates high and blue color low importance. The central location is the original target location, hence the central description identical to Fig. 3A.
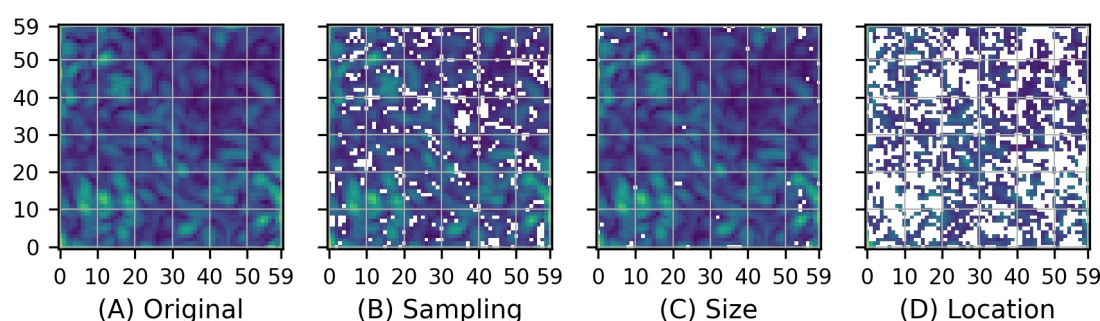
449 this checkerboard pattern is indeed a known artifact of strided convolutions and max-pooling layers used in
450 CNN2 (Odena et al., 2016).

451     Figure 6 shows the same as Fig. 3 and 5 but for the logistic regression model. For the sampling variant
452 tasks, large parts of the original description $\vec{d}$ are masked. This indicates that these parts likely reflect
453 spurious relations. For the size variant task, on the other side, only few pixels are masked. Lastly, for the
454 location variant tasks, nearly all pixels are masked. This indicates that the original description $\vec{d}$ shown in
455 Fig. 6A likely reflects spurious relations only.

456     For this task, we know that the physical importance of a pixel averaged over a long time period decreases
457 with the pixel's distance to the central target patch. Further, due to the predominantly westerly winds, the

**Figure 5.** Same as Fig. 3 but for CNN2.



**Figure 6.** Same as Fig. 3 but for the logistic regression model.

458 average physical importance of pixels is slightly shifted to the west. Given this knowledge, we can confirm
459 that the variant approach successfully identified all pixels in Fig. 3A, 5A, and 6A which reflect spurious
460 relations. Note that the sampling approach alone (see Fig. 3B, 5B, and 6B), which is the commonly applied
461 method, is not sufficient to identify all pixels reflecting spurious relations.

462 Note further that the examples emphasize once again the following: even if parts of a description are not
463 indicated as spurious by any considered variant task, we cannot conclude that they reflect causal relations.
464 Imagine, for instance, that we had only considered the size variant task. For this variant task and the
465 logistic regression model, only a small number of pixels is masked although Fig. 6A seems to exclusively
466 reflect spurious relations. Hence, variant tasks can only indicate parts of an original description as likely
467 reflecting spurious relations and do not allow for any direct inference about other parts of the description.
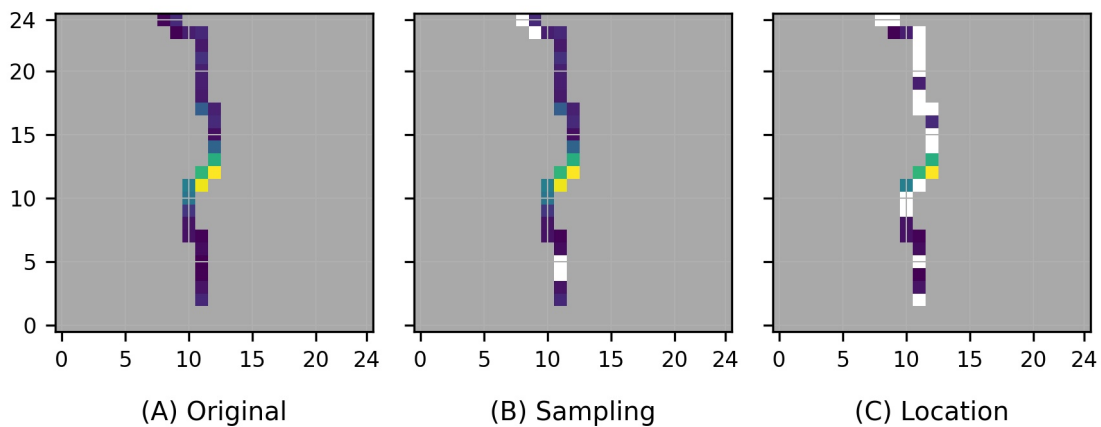468 Nevertheless, this can be useful already.

## 3.2 Task 2 – Water level prediction

469

470 Figure 7a shows the description $\vec{d}$ of the function that the MLP learned when it was trained on the
471 original water level prediction task. Remember that the considered description is a measure of the average
472 importance of each pixel in the considered river segment for the predictions of the model. Our objective
473 is to apply the variant approach to identify parts of the description that reflect spurious relations. To that
474 purpose we computed the variant descriptions $\vec{d^{v_i}}$ for all sampling and location variant tasks, and masked
475 all pixels of $\vec{d}$ for which the relative distance between the original description $\vec{d}$ and one of the (shifted)
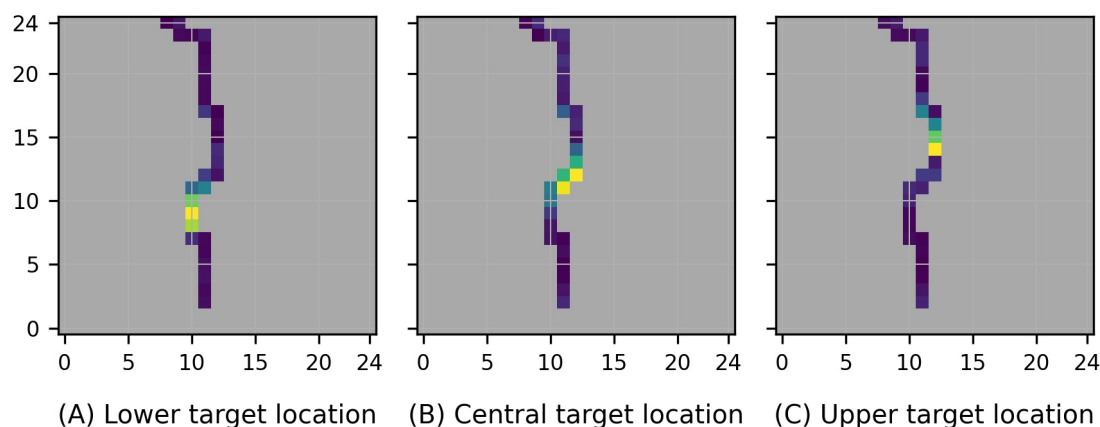
476 variant descriptions exceeds the threshold of $t = 0.5$. The resulting masked versions of Fig. 7A are shown
477 in Fig. 7B and C.

478     For this task, we know that the development of the water level at the target location depends only on
479 the water level closely upstream and downstream. Hence, Fig. 7A is (apart from the moderately high
480 importance of pixel (11,17)) close to our understanding of the physical importance of the considered pixels.
481 Nevertheless, especially in Fig. 7C, many of the pixels further upstream and downstream of the target
482 location are masked, i.e. the variant approach indicates (mistakenly) that the low feature importance of
483 these pixels likely reflects spurious relations. We suspect that this happened because we considered relative
484 rather than absolute distances between original and variant descriptions (see Eq. 1), which can cause two
485 small values to have a large distance which in turn causes the corresponding pixel to be mistakenly masked
486 as spurious. Apart from pixels with low feature importance, also pixel (11,11) closely upstream of the
487 original target location seems to be mistakenly masked as spurious in Fig. 7C. We suspect that this is due
488 to our assumption that causal relations are shifted along the river by the exact same number of pixels as the
489 target location is. While this assumption enables us to simply consider pixelwise relative distances between
490 original description $\vec{d}$ and shifted variant descriptions $\vec{d^{v_i}}$ (see Sect. Methods), it might be overly simplified
491 as for example the flow velocity at different locations in the river might differ, and the river might cross
492 some pixels diagonally and others straight.

493     Here, a visual assessment of the individual variant descriptions seems to be superior to the formal
494 evaluation of distances performed for Fig. 7C because it allows a softer comparison between original and
495 variant descriptions $\vec{d}$ and $\vec{d^{v_i}}$. Indeed, upon visual assessment of the location variant descriptions depicted
496 in Fig. 8, and with the assumption in mind that causal relations *approximately* reflect the shift of the target
497 location, the only pixel in Fig. 7A that we would mark as potentially reflecting spurious relations, is pixel
498 (11,17).



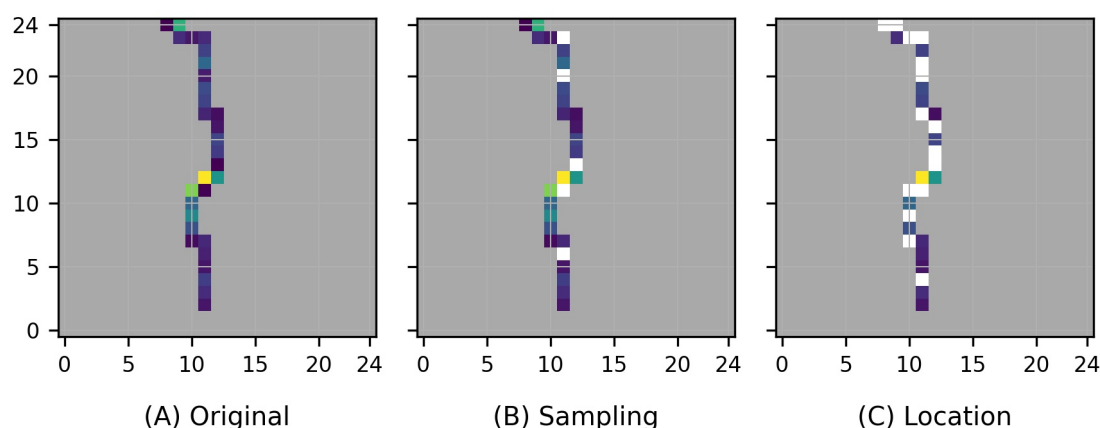(A) Original      (B) Sampling      (C) Location

**Figure 7.** **(A)** Description $\vec{d}$ of the function that the MLP learned when it was trained on the original water
level prediction task (see Fig. 1B). The description is a measure of the average importance of each pixel in
the considered river segment for the predictions of the model. Yellow color indicates high and blue color
low importance. **(B-C)** As (A), but pixels for which the relative distance between the original description $\vec{d}$
and one of the sampling and (shifted) location variant descriptions $\vec{d^{v_i}}$, respectively, exceeds the threshold
of $t = 0.5$, are masked. Gray marks pixels outside the considered river segment.
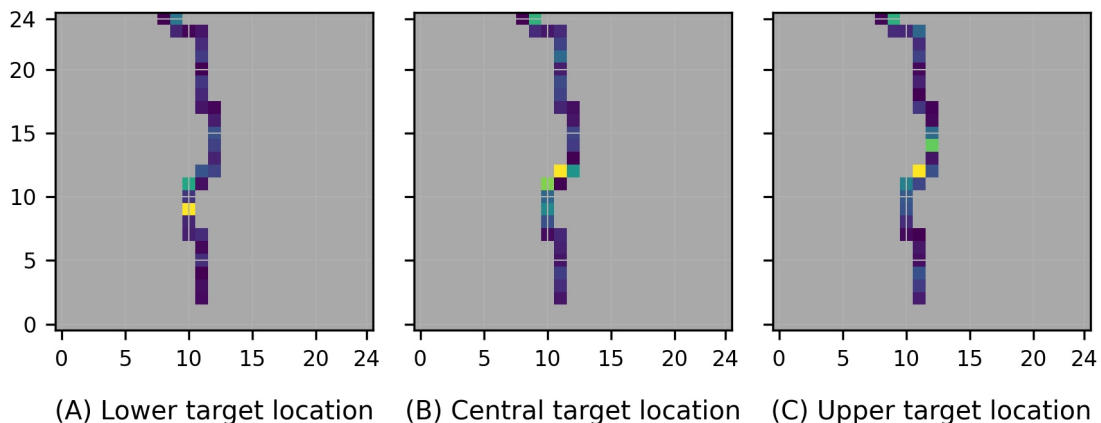
**Figure 8.** Descriptions of the functions that separate instances of the MLP learned when trained for the different target locations (from left to right the target location is at (10,10), (12,12), (12,15), see Fig. 2B). Note that panel B shows the same as Fig. 7A. Gray marks pixels outside the considered river segment.

499    Figures 9 and 10 show the same as Figs. 7 and 8 but for the linear regression model. In this case, the
500    formal evaluation of distances between original and location variant descriptions performed for Fig. 9C
501    indicates that Fig. 9A reflects spurious relations at nearly all pixels except from the target location and the
502    neighboring pixel upstream. In this case, the formal evaluation agrees well with the visual assessment of
503    the location variant descriptions depicted in Fig. 10. Indeed, visual assessment of Fig. 10 also indicates that
504    the neighboring pixel upstream of the target location and maybe the target location itself are the only two
505    pixels for which the assigned importance approximately reflects the shift of the target location between
506    Fig. 10A, B and C.



**Figure 9.** Same as Fig. 7 but for the linear regression model.

(A) Lower target location    (B) Central target location    (C) Upper target location

**Figure 10.** Same as Fig. 8 but for the linear regression model.

## 4   CONCLUSIONS

Given a description $\vec{d} \in \mathbb{R}^d$ of the function that a statistical model learned during a training phase, we proposed a variant approach for the identification of parts of $\vec{d}$ that reflect spurious relations. We successfully demonstrated the approach and its superiority over pure sampling approaches with two illustrative hydrometeorological predictions tasks, various statistical models and illustrative descriptions. For the rain prediction task, where we assumed causal relations to remain stable between original and variant tasks, the formal evaluation of distances between original and variant descriptions enabled us to correctly identify all spurious relations that the statistical models learned. For the water level prediction task, where formally specifying the assumed variation of causal relations was more involved, we found the formal evaluation of distances to be of limited use. However, visual assessment enabled us again to correctly identify all spurious relations that the statistical models learned.

In this work, we considered simplified tasks and global descriptions of the learned functions to be able to decide whether parts of the descriptions that the variant approach identifies as spurious do indeed reflect spurious relations. This was necessary to evaluate the variant approach. However, we expect the approach to be beneficial for a wide range of more complex prediction tasks. Naming two possible applications outside the geosciences, it might be used to identify spurious relations reflected in (local) descriptions of functions that DL models trained on electroencephalography (EEG) data (Sturm et al., 2016) learned by comparing them to variant descriptions obtained for variant models trained for different (groups of) patients; or to automatically detect spurious relations reflected in (local) descriptions of functions that a DL model trained on a common image data set learned (Lapuschkin et al., 2019) by automatically comparing them to variant descriptions for variant models trained on different image data sets. Applications of the variant approach to more complex prediction tasks in the geosciences and beyond, and to local descriptions of the learned functions, are planned in future.

A challenge when applying the proposed variant approach may be to define variant tasks beyond random sampling of the training data. However, a data set is often composed of different sources constituting in themselves variants. Further, the modification of the rain prediction task, where we were able to identify

532 parts of the original description as spurious by merely changing the size of the input region, indicates that
533 even small modifications of the original prediction task can be useful.

534     Apart from the variant approach, which considers a fixed statistical model and modifications of an
535 original prediction task, another approach for identifying spurious relations that a considered statistical
536 model learned might be to compare the relations between input and target variables that different models
537 learn when trained on the (fixed) prediction task. In such an approach, the degree of variation between
538 models may differ from varying configurations in Monte-Carlo dropout, over random seeds for the weight
539 initialization of otherwise identical models to completely different statistical models. Formalization and
540 evaluation of this approach is out of scope of this work.

## 5   CONFLICT OF INTEREST STATEMENT

541 The authors declare that the research was conducted in the absence of any commercial or financial
542 relationships that could be construed as a potential conflict of interest.

## 6   AUTHOR CONTRIBUTIONS

543 T.T. and S.K. designed the experiments. T.T. conducted the experiments, analyzed the results and prepared
544 the manuscript with contributions from S.K. and J.G..

## 7   ACKNOWLEDGMENTS

## 8   DATA AND CODE AVAILABILITY STATEMENT

555 The original contributions presented in the study are publicly available. The data and code can be found
556 here `https://datapub.fz-juelich.de/slts/t_tesch/`.

## REFERENCES

557 Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K., and Samek, W. (2015). On pixel-wise
558     explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLOS ONE* 10,
559     e0130140. doi:10.1371/journal.pone.0130140
560 De Bin, R., Janitza, S., Sauerbrei, W., and Boulesteix, A.-L. (2015). Subsampling versus bootstrapping in
561     resampling-based model selection for multivariable regression. *Biometrics* 72, 272–280. doi:10.1111/
562     biom.12381

Furusho-Percot, C., Goergen, K., Hartick, C., Kulkarni, K., Keune, J., and Kollet, S. (2019). Pan-european groundwater to atmosphere terrestrial systems climatology from a physically consistent simulation. *Sci. Data* 6, 320. doi:10.1038/s41597-019-0328-7

Gagne II, D. J., Haupt, S. E., Nychka, D. W., and Thompson, G. (2019). Interpretable deep learning for spatial analysis of severe hailstorms. *Mon. Wea. Rev.* 147, 2827–2845. doi:10.1175/mwr-d-18-0316.1

Gasper, F., Goergen, K., Shrestha, P., Sulis, M., Rihani, J., Geimer, M., et al. (2014). Implementation and scaling of the fully coupled terrestrial systems modeling platform (terrsysmp v1.0) in a massively parallel supercomputing environment – a case study on juqueen (ibm blue gene/q). *Geosci. Model Dev.* 7, 2531–2543. doi:10.5194/gmd-7-2531-2014

Gilpin, L., Bau, D., Yuan, B., Bajwa, A., Specter, M., and Kagal, L. (2018). Explaining explanations: An overview of interpretability of machine learning. In *2018 IEEE 5th International Conference on Data Science and Advanced Analytics (DSAA)* (IEEE), 80–89. doi:10.1109/dsaa.2018.00018

Goodfellow, I., Shlens, J., and Szegedy, C. (2015). Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*

Guo, R., Cheng, L., Li, J., Hahn, P. R., and Liu, H. (2020). A survey of learning causality with data: Problems and methods. *ACM Comput. Surv.* 53. doi:10.1145/3397269

Ham, Y., Kim, J., and Luo, J. (2019). Deep learning for multi-year ENSO forecasts. *Nature* 573, 568–572. doi:10.1038/s41586-019-1559-7

He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 770–778. doi:10.1109/CVPR.2016.90

Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *32nd International Conference on Machine Learning, PMLR*. vol. 37, 448–456

Lakkaraju, H., Arsov, N., and Bastani, O. (2020). Robust and stable black box explanations. In *37th International Conference on Machine Learning, PMLR*. vol. 119

Lapuschkin, S., Wäldchen, S., Binder, A., Montavon, G., Samek, W., and Müller, K. (2019). Unmasking clever hans predictors and assessing what machines really learn. *Nat. Commun.* 10, 1096. doi:10.1038/s41467-019-08987-4

Larraondo, P. R., Renzullo, L. J., Inza, I., and Lozano, J. A. (2019). A data-driven approach to precipitation parameterizations using convolutional encoder-decoder neural networks

LeCun, Y., Bottou, L., Orr, G., and Müller, K. (2012). *Efficient BackProp* (Berlin, Heidelberg: Springer Berlin Heidelberg). 9–48. doi:10.1007/978-3-642-35289-8_3

McGovern, A., Lagerquist, R., Gagne II, D. J., Jergensen, G. E., Elmore, K. L., Homeyer, C. R., et al. (2019). Making the black box more transparent: Understanding the physical implications of machine learning. *Bull. Amer. Meteor. Soc.* 100, 2175–2199. doi:10.1175/bams-d-18-0195.1

Molnar, C. (2019). *Interpretable Machine Learning*

Montavon, G., Samek, W., and Müller, K. (2018). Methods for interpreting and understanding deep neural networks. *Digit. Signal Process.* 73, 1–15. doi:10.1016/j.dsp.2017.10.011

Odena, A., Dumoulin, V., and Olah, C. (2016). Deconvolution and checkerboard artifacts. *Distill* doi:10.23915/distill.00003

Pan, B., Hsu, K., AghaKouchak, A., and Sorooshian, S. (2019). Improving precipitation estimation using convolutional neural network. *Water Resour. Res.* 55, 2301–2321. doi:10.1029/2018wr024090

Pan, S. J. and Yang, Q. (2010). A survey on transfer learning. *IEEE Trans. Knowl. Data Eng.* 22, 1345–1359. doi:10.1109/tkde.2009.191

Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., et al. (2019). Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*

*32*, eds. H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. Fox, and R. Garnett (Curran Associates, Inc.). 8026–8037

Pearl, J. (2009). Causal inference in statistics: An overview. *Statistics Surveys* 3. doi:10.1214/09-ss057

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., et al. (2011). Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* 12, 2825–2830

Peters, J., Bühlmann, P., and Meinshausen, N. (2016). Causal inference by using invariant prediction: identification and confidence intervals. *J. R. Stat. Soc.: Series B (Statistical Methodology)* 78, 947–1012. doi:10.1111/rssb.12167

Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). "Why should i trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (New York, NY, USA: Association for Computing Machinery), KDD '16, 1135–1144. doi:10.1145/2939672.2939778

Roscher, R., Bohn, B., Duarte, M. F., and Garcke, J. (2020). Explainable machine learning for scientific insights and discoveries. *IEEE Access* 8, 42200–42216. doi:10.1109/access.2020.2976199

Samek, W., Montavon, G., Lapuschkin, S., Anders, C. J., and Müller, K. R. (2021). Explaining Deep Neural Networks and Beyond: A Review of Methods and Applications. *Proceedings of the IEEE* 109, 247–278. doi:10.1109/JPROC.2021.3060483

Schramowski, P., Stammer, W., Teso, S., Brugger, A., Herbert, F., Shao, X., et al. (2020). Making deep neural networks right for the right scientific reasons by interacting with their explanations. *Nat. Mach. Intell.* 2, 476–486. doi:10.1038/s42256-020-0212-3

Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. (2017). Grad-CAM: Visual explanations from deep networks via gradient-based localization. In *2017 IEEE International Conference on Computer Vision (ICCV)* (IEEE), 618–626. doi:10.1109/iccv.2017.74

Shimodaira, H. (2000). Improving predictive inference under covariate shift by weighting the log-likelihood function. *J. Stat. Plan. Inference* 90, 227 – 244. doi:10.1016/s0378-3758(00)00115-4

Shrestha, P., Sulis, M., Masbou, M., Kollet, S., and Simmer, C. (2014). A scale-consistent terrestrial systems modeling platform based on cosmo, clm, and parflow. *Mon. Wea. Rev.* 142, 3466–3483. doi:10.1175/MWR-D-14-00029.1

Simonyan, K., Vedaldi, A., and Zisserman, A. (2013). Deep inside convolutional networks: Visualising image classification models and saliency maps.

Sinha, A., Namkoong, H., and Duchi, J. C. (2018). Certifiable distributional robustness with principled adversarial training. In *International Conference on Learning Representations*

Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* 15, 1929–1958

Sturm, I., Lapuschkin, S., Samek, W., and Müller, K. (2016). Interpretable deep neural networks for single-trial EEG classification. *J. Neurosci. Methods* 274, 141–145. doi:10.1016/j.jneumeth.2016.10.008

Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., et al. (2014). Intriguing properties of neural networks. In *International Conference on Learning Representations*

Toms, B. A., Barnes, E. A., and Ebert-Uphoff, I. (2020). Physically interpretable neural networks for the geosciences: Applications to earth system variability. *Journal of Advances in Modeling Earth Systems* 12, e2019MS002002. doi:10.1029/2019ms002002

Zhang, Q. and Zhu, S. (2018). Visual interpretability for deep learning: a survey. *Front. Inf. Technol. Electron. Eng* 19, 27–39. doi:10.1631/fitee.1700808