

# Deep Learning-based Parameter Transfer in Meteorological Data

Fatemeh Farokhmanesh,<sup>a</sup> Kevin Höhlein,<sup>a</sup> and Rüdiger Westermann<sup>a</sup>

<sup>a</sup>*Department of Informatics, Technical University of Munich, Germany*

<sup>4</sup> Corresponding author: Fatemeh Farokhmanesh, fatemeh.farokhmanesh@tum.de

5 ABSTRACT: Numerical simulations in earth-system sciences consider a multitude of physical  
6 parameters in space and time, leading to severe I/O bandwidth requirements and challenges in  
7 subsequent data analysis tasks. Deep-learning based identification of redundant parameters and  
8 prediction of those from other parameters, i.e. Variable-to-Variable (V2V) transfer, has been  
9 proposed as an approach to lessening the bandwidth requirements and streamlining subsequent  
10 data analysis. In this paper, we examine the applicability of V2V to meteorological reanalysis  
11 data. We find that redundancies within pairs of parameter fields are limited, which hinders  
12 application of the original V2V algorithm. Therefore, we assess the predictive strength of reanalysis  
13 parameters by analyzing the learning behavior of V2V reconstruction networks in an ablation  
14 study. We demonstrate that efficient V2V transfer becomes possible when considering groups of  
15 parameter fields for transfer, and propose an algorithm to implement this. We investigate further  
16 whether the neural networks trained in the V2V process can yield insightful representations of  
17 recurring patterns in the data. The interpretability of these representations is assessed via layer-  
18 wise relevance propagation that highlights field areas and parameters of high importance for the  
19 reconstruction model. Applied to reanalysis data, this allows uncovering mutual relationships  
20 between landscape orography and different regional weather situations. We see our approach as  
21 an effective means to reduce bandwidth requirements in numerical weather simulations, which can  
22 be used on top of conventional data compression schemes. The proposed identification of multi-  
23 parameter features can spawn further research on the importance of regional weather situations for  
24 parameter prediction, also in other kinds of simulation data.

25 **1. Introduction**

26 The rapid increase in available computing power has enabled a broad adoption of simulation-  
27 based research methodologies in earth-system sciences. Numerical simulations of spatio-temporal  
28 dynamical systems consider a multitude of physical parameters and are carried out at high resolution  
29 in space and time. Especially in meteorology and climate modelling, also numerical ensemble  
30 simulations are carried out with varying magnitudes of initial condition uncertainty, to account  
31 for the uncertainty in the representation of certain physical processes. Simulations are performed  
32 routinely by weather centers worldwide, and in research we see increasing use of unique super-  
33 ensembles consisting of hundreds and even thousands of members (Necker et al. 2020).

34 In classical workflow scenarios, simulations are run on large-scale computing facilities and data  
35 are streamed to and stored on external file systems for archiving and subsequent analysis. However,  
36 the volume of generated data has reached an order of magnitude where the speed of data transfer  
37 between computing device and file system – so-called I/O operations – imposes a major bottleneck.  
38 For instance, during the 2010's, the ability to compute (and thus generate data) increased on  
39 supercomputers much faster than the ability to store and load data, roughly about two orders of  
40 magnitude over this decade. I/O performance, in contrast, only increased one order of magnitude.

41 The divergence between compute and I/O renders the classical simulation workflow increasingly  
42 problematic, and requires – in addition to using classical data compression techniques – to avoid  
43 streaming or even simulating redundant data that can be recovered from the generated results. Such  
44 an approach requires less I/O bandwidth, also for bringing the data to the compression stage, and  
45 can improve the time it requires to bring the data into a format suitable for the used compressor.

46 Within this line of research, deep-learning-based Variable-to-Variable (V2V) transfer has been  
47 proposed recently by Han et al. (2021b) for optimizing information transfer in situations were  
48 spatio-temporal multi-parameter simulations can be carried out in far less time than it requires to  
49 store the data on a file system. V2V considers each simulated parameter as a separate entity and  
50 proposes an algorithm to identify groups of similar parameters and one representative member  
51 from which the other parameters in this group can be inferred.

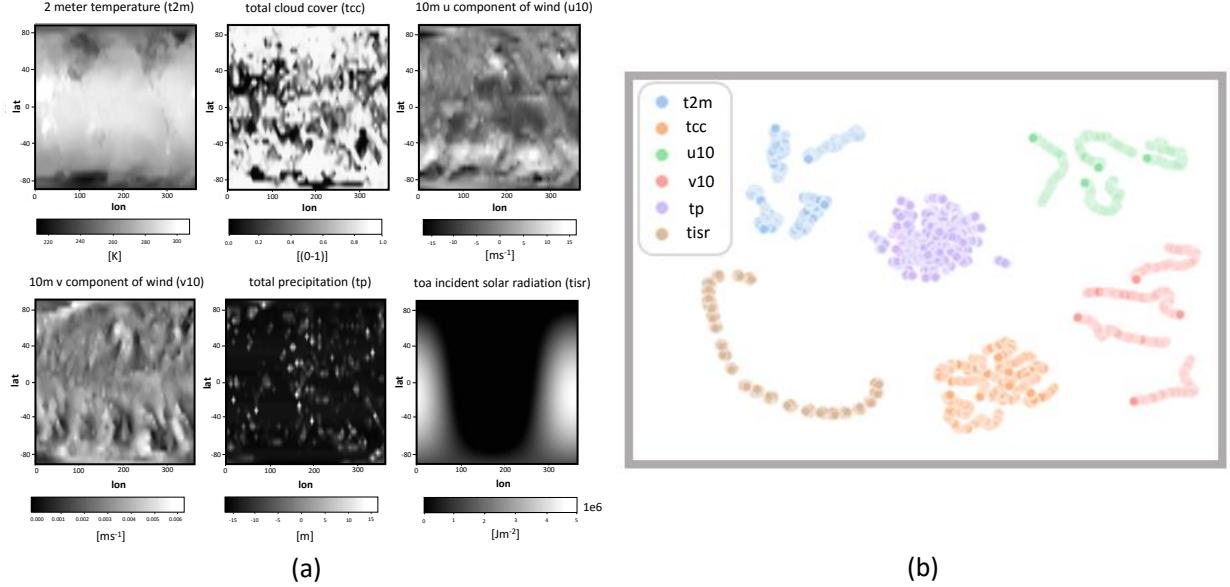
52 V2V represents all simulated parameter fields in a common feature space (the so-called latent-  
53 space) that is learnt by a convolutional neural network (CNN), and identifies subsets of similar  
54 parameters in this space. For each subset the most representative member is determined, and

another suitably trained network then learns to reconstruct all other member in one subset from the representative one. The bandwidth requirements for storing the multi-parameter data is reduced to the bandwidth required for storing the representative parameter fields and the weight parameterization of the reconstruction networks.

## 59 *Contribution*

In this work, we assess the applicability of V2V transfer to earth-system related data, and identify shortcomings of the proposed methodology. In consideration of these findings, we propose an improved approach, in which a single CNN is trained on meteorological archives to learn general relationships between subsets of parameters, thereby focusing on subsets of parameter fields with vastly different characteristics and variation in 2D space and time. We demonstrate the capabilities of the proposed approach using two exemplary datasets, representing different modalities of meteorological data. We consider a global reanalysis dataset, which is taken from the WeatherBench (WB) benchmark suite (Rasp et al. 2020), and study the performance of the proposed approach on data on large spatial scales. Furthermore, we apply the proposed method to an ensemble forecast dataset, which was generated by Necker et al. (2020) to study the sampling accuracy of spatio-temporal correlation patterns in convective-scale forecast ensembles. Fig. 1, as well as Fig. B1, demonstrate significant structural variability between the single-parameter fields in each dataset. Fig. 1 (a) shows snapshots of the parameter fields at a particular point in time. Fig. 1 (b) shows a two-dimensional embedding (computed with t-SNE, van der Maaten and Hinton 2008) of the V2V latent-space representations of these fields at different time points, which are used to search for parameter similarities. Fig. B1 displays the same overview for the convective-scale ensemble (CSEns) dataset. It can be seen that visible clusters involve only a single parameter, and clusters of distinct parameters are situated at roughly the same distance from one another. The lack of similarities between pairs of parameters prohibits the use of the original V2V algorithm.

Based on these observations, we propose a different strategy for V2V transfer, which considers the expressiveness of subsets of multiple parameters instead of transferring only between pairs of them. To identify these subsets, a CNN-based model architecture is trained multiple times on multi-parameter fields with varying parameter subsets removed from the input data. In an ablation study, we then shed light on the prediction skills of the different models, depending on



79 FIG. 1. Different parameter fields in the ERA5 reanalysis dataset. a) Gray-scale visualizations of the parameter  
 80 fields at a particular time. b) t-SNE projections of latent-space features of the parameter fields (different  
 81 parameters indicated by colors) at different times (note that projections for different initializations of t-SNE yield  
 82 similar groupings).

88 the parameter subsets that are used as source. For  $n \in \mathbb{N}$  originally available parameter fields, the  
 89 networks are designed to learn mappings from  $m (< n)$  parameter fields (the input) to predict the  
 90 remaining  $n - m$  ones (the output). Doing so, the networks encode the inputs into a compact latent-  
 91 space representation, in which relationships between the  $m$  input fields and the  $n - m$  output fields  
 92 are encoded. The networks are trained via standard backpropagation with a loss function, which  
 93 measures the reconstruction accuracy for all  $n - m$  parameters in common. We demonstrate, using  
 94 numerical and visualization-based quality metrics, that the networks efficiently learn to reconstruct  
 95 the unseen parameters, thereby not overfitting to the provided training samples but generalizing to  
 96 simulation snapshots that haven't been seen before.

97 Due to the high computational complexity of training all  $\binom{n}{m}$  different networks for reconstruct-  
 98 ing  $n - m$  fields from  $m$  given fields, we propose a computationally less involved strategy and  
 99 demonstrate its effectiveness for selecting the most representative members. Conceptually, this  
 100 strategy builds upon removing iteratively those parameters that are most difficult to predict from the  
 101 remaining parameters, simultaneously avoiding to keep redundant fields in the input. Furthermore,

102 the networks' training behaviours are monitored and the convergence rates at early training stages  
103 after few epochs of learning are used as indicators of the difficulty of parameter transfer. For  
104 instance, for one of our use cases comprising nine different parameters and selecting four of them  
105 to predict the remaining five parameters, training requires only roughly six hours on a low-size  
106 deep learning cluster with six mid-size GPUs, equipped with 11 GB of graphics memory, each.

107 Beyond considering the networks purely as black-box models, we further try to gain insight into  
108 the multi-parameter relationships that are learned by the models. For this purpose, we employ  
109 an adapted version of layer-wise relevance propagation (LRP, Bach et al. 2015), a method for  
110 highlighting input areas and parameters, which are important for the reasoning process of the  
111 models. This offers the opportunity to uncover feature patterns in multi-parameter space, i.e.  
112 reoccurring parameter combinations, which are recognized as important for the network to achieve  
113 high accuracy, and analyse their correspondence to certain weather situations.

114 The remainder of this paper is structured as follows. In section 2, we review related work.  
115 In section 3, we summarize the original V2V algorithm and highlight algorithmic shortcomings.  
116 Building up on these findings, we present our extended V2V approach in section 4. We introduce  
117 the example datasets used for subsequent experiments in section 5 and describe the network  
118 architectures used for the experiments in section 6. The ablation study, as well as the LRP analysis  
119 of the models, are carried out in section 7. We conclude the paper in section 8.

## 120 2. Related work

121 In recent years, machine learning with powerful deep-learning architectures has found applica-  
122 tions various fields of climate science and meteorology (Reichstein et al. 2019). Many of the  
123 possible applications exploit the efficiency and flexibility of CNN architectures when applied to  
124 inference tasks involving grid-structured data.

125 Related to our approach are so-called super-resolution and downscaling techniques, which re-  
126 construct high-resolution parameter fields from corresponding low-resolution versions. In contrast  
127 to V2V approaches, the information transfer occurs between representations of the same parameter  
128 with different spatial resolutions. Some of these approaches can be used for data compression, in  
129 principle. Such methods operate by first sub-sampling the parameter fields and subsequently re-  
130 constructing the initial fields from the sub-sampled versions. For example, Rodrigues et al. (2018)

proposed a supervised convolutional neural network that interpolates a low-resolution weather data into a high-resolution output. Pouliot et al. (2018) introduced deep learning-based enhancement in landsat super-resolution. Cheng et al. (2020) proposed a method that converts low-resolution climate data to high-resolution climate forecasts using Laplacian pyramid super-resolution networks. Downscaling approaches, in contrast, aim at predicting additional high-resolution details from low-resolution parameter fields, without assuming prior knowledge of the original high-resolution data. For instance, Höhlein et al. (2020) and Serifi et al. (2021) train on a small set of paired low- and high-resolution simulation pairs to circumvent generating expensive high-resolution simulations at inference time at all. These techniques, even they establish relationships between low- and high-resolution fields, are not motivated by the idea to compress the data.

Super-resolution of scientific data has also been investigated from the perspective of scientific visualization, since high-resolution simulations in meteorology make analysis of these datasets challenging. Early works on data super-resolution demonstrate the capabilities of neural networks to learn upscaling a low-resolution version of the data to the initial high-resolution dataset. Upscaling is performed in the spatial domain (Zhou et al. 2017; Han and Wang 2020; Guo et al. 2020), in the temporal domain (Han and Wang 2019), and in the spatio-temporal domain (Han et al. 2021a). Underlying these works is the goal to avoid storing the high-resolution datasets and, thus, reduce bandwidth and memory requirements. By training networks to infer the full image from a low-resolution image of an iso-surface, Weiss et al. (2019) demonstrate improved rendering frame rates. In recent work by Weiss et al. (2020) a convolutional neural network learns to adaptively place image samples and reconstruct the full image from the generated unstructured set of samples.

To work in situations with severe I/O bottleneck, Sato et al. (2019) introduced an in-situ processing strategy for visualization and post-processing of high-resolution meteorological data. Röber and Engels (2019) analysed in-situ data processing approaches in climate science. Helbig et al. (2015) proposed a visualization workflow where the first stage is a data abstraction layer that downsamples the data spatially and temporally. Toderici et al. (2017) proposed an image compression method using a recurrent neural network by saving the compressed latent space produced by the network instead of the high-resolution data.

A different approach for data compression has been introduced by Han et al. (2021b) for multi-parameter data, by training networks to infer certain parameters from others, and, thus, to avoid

161 storing these parameters. Conceptually, this work builds upon the notion of information transfer  
162 between scalar fields (Wang et al. 2011) to derive transferable parameter pairs.

### 163 3. Deep-learning-based variable-to-variable transfer

#### 164 a. V2V Algorithm

165 The overall goal of V2V lies in identifying those variables in a multi-parameter dataset that can  
166 be redundantly reconstructed given other parameters from the same dataset. Han et al. (2021b)  
167 subdivide the V2V process into 3 conceptually distinct stages: feature learning, translation graph  
168 construction and variable translation. In the feature learning stage, a CNN model architecture  
169 such as UNet (Ronneberger et al. 2015) is trained in an auto-encoder-like setting to encode and  
170 reconstruct snapshots of parameter fields. The same model is shared between all parameters, such  
171 that the internal hidden variables of the model, referred to as the latent-space representation, are  
172 informed about similarities and dissimilarities between different parameters. After training, all  
173 parameter snapshots are mapped into the latent space where clusters of similar parameters are  
174 detected. Han et al. (2021b) propose to find clusters through visual examination of the latent-space  
175 features. To visualize the features, they apply a non-linear dimension reduction algorithm, called  
176 t-distributed stochastic neighbor embedding (t-SNE, van der Maaten and Hinton 2008). Parameters  
177 in a common cluster become a transferable variable group.

178 In the translation-graph construction stage, the Kullback-Leibler divergence is used to estimate a  
179 measure of so-called *transferable difficulty* for pairs of parameters inside a transferable parameter  
180 group. Parameter pairs are considered for transfer, if the Euclidean distance between their respective  
181 latent-feature representations is smaller than a predefined distance threshold. Parameters in different  
182 transferable parameter groups are not considered for transfer.

183 A directed transfer graph is then constructed by chaining transferable pairs according to a  
184 minimum discrepancy criterion. Finally, in the variable transfer stage, CNNs are trained to learn  
185 the transfer mapping according to the translation graph.

#### 186 b. Shortcomings of V2V

187 V2V reduces the search for transferable variables to pairs of similar variables based on a distance  
188 threshold criterion and visual analysis of a t-SNE projection. This leads to a number of shortcom-

ings regarding expressivity of the identified parameter relations and reproducibility of the results.  
In a first place, V2V does not assess true predictability of one variable based on another, but uses an empirically determined approximative criterion, which might overlook potentially valuable relationships when the threshold value is not set optimally.

The same holds true for subsets of dissimilar variables from which another dissimilar variable might be predictable. Since V2V is constrained to pairwise similarities, higher-order similarities between multiple parameters cannot be considered. Furthermore, V2V cannot always decide unambiguously the source and target fields. This decision is based on a so-called translation graph where the nodes correspond to the parameter fields and directed edges indicate transferability from a source to a target variable. In this graph, however, cycles can occur. This happens, because the *transferable difficulty* is based on Kullback-Leibler divergences, which are not a metric and don't satisfy the triangle inequality in general. When only pairwise transferabilities are considered, this can result in the selection of all parameters in a set of similar parameters as sources and targets of one another, respectively. Concerning reproducibility, the use of t-SNE brings a number of potential problems. t-SNE is a nonlinear dimension reduction algorithm and the low-dimensional projections obtained from it are known to depend non-trivially on a number of hyper-parameters and weight initializations. Depending on the specific realization of hyper-parameters and initialization, projections may differ largely. This, in turn, affects the identification of transferable parameter groups, which are used to set important constraints on parameter comparisons in the translation graph construction stage. Besides this, manual interaction with the running simulation, i.e. parameter selection at runtime, may not be a favorable option in operational high-performance computing scenarios.

## 4. Method

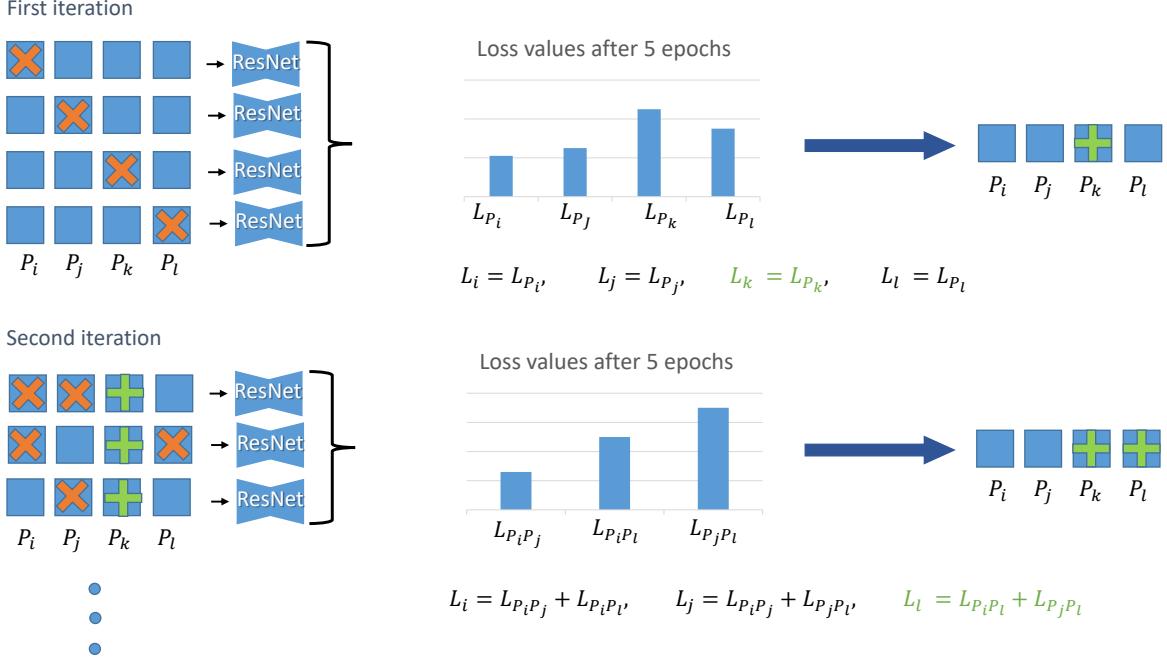
To overcome the shortcomings of V2V, we propose an alternative parameter selection procedure. It replaces the single-parameter auto-encoding CNN and subsequent clustering and pairwise similarity search with multi-parameter CNNs and loss-based tracking of the learning progress. Given a multi-variate time-varying dataset with  $n \in \mathbb{N}$  parameters and  $T \in \mathbb{N}$  timesteps, the user selects the number  $m$  of input fields from which the remaining  $m - n$  parameter fields are predicted. The most straight forward – yet computationally demanding – approach is to launch  $\binom{n}{m}$  training runs

218 for the different parameter configurations, and select the network with the lowest loss. By starting  
219 from  $n - 1$  inputs and one output and proceeding iteratively with decreasing number of inputs, the  
220 procedure can also be made dependent on a predefined loss threshold, i.e., by launching for all  
221  $k$ ,  $1 \leq k \leq m$ , a batch of  $\binom{n}{k}$  training runs and stopping once the minimum loss exceeds a given  
222 threshold. The parameter configuration for which the minimum loss is achieved is then selected  
223 for variable transfer. In our current implementation we consider only the user-specified number of  
224 input parameters.

225 Since for large values of  $n$ , the described procedure requires training too many networks, a  
226 computationally less expensive alternative needs to be developed. A straight forward approach is  
227 to train only  $n$  networks, where each network predicts one single parameter from the remaining  
228  $n - 1$  parameters, and use the networks' losses as indicators of how difficult the prediction of a  
229 parameter is. If a loss value is high, predictability is low, and, thus, the parameter predicted with  
230 the highest loss is fixed as one of the input parameters. Then, of all trained networks where this  
231 parameter is already contained in the input set, the one with the largest loss is selected, and the  
232 variable that is predicted is added to the input. This procedure is repeated until the specified  
233 number of inputs is reached.

234 Interestingly, in our experiments this approach finds exactly the same input and output sets as  
235 the exhaustive training procedure, yet it requires training only  $n$  networks. In general, however,  
236 since only the difficulties of predicting each parameter individually from all other parameters are  
237 considered, the prediction strength of a reduced parameter set can be overlooked or overestimated.  
238 For instance, consider a subset of similar parameters, each of which can be predicted at high accu-  
239 racy from the other parameters in this subset. Consequently, the loss values of the corresponding  
240 networks will be low, so that it becomes unlikely that one of these parameters will be selected as  
241 input. Then, however, the input may solely comprise parameters from which the ones in the subset  
242 cannot be well predicted.

243 To address these shortcoming, we propose a strategy with lower computational complexity than  
244 the first strategy, and which differs from the second strategy in that it considers subsets of parameters  
245 in both the inputs and predicted outputs. As in the previous strategy,  $n$  networks are trained initially,  
246 which each network predicting one single parameter from the remaining  $n - 1$  parameters, and the  
247 parameter predicted by the network with the highest loss is fixed in the input. In the next iteration,



243 FIG. 2. Method overview. In the  $i$ -th iteration, the same Resnet is trained multiple times using different  
 244 combinations of  $n - i$  input and  $i$  output parameters. Orange crosses indicate the output parameters, green pluses  
 245 indicate the parameters that are fixed in the input set, blue squares indicate the remaining parameters in the input  
 246 set. In each iteration, for each free parameter a loss is computed by adding the losses of those networks in which  
 247 the parameter is in the output set. The parameter with the maximum loss is fixed in the input set. For the WB  
 248 dataset, orography is used in every input, but is not predicted.

254 all networks with  $n - 2$  inputs (including the fixed parameter) and two outputs are trained. For all  
 255 but the fixed parameter, the overall loss is computed by adding the losses of all networks where this  
 256 parameter is in the predicted set. The parameter with the highest loss is fixed, and the procedure  
 257 moves on with  $n - 3$  inputs, and three outputs now containing the two fixed parameters (see Fig. 2  
 258 for a graphical overview of the proposed approach). This strategy, in case of a similar subset of  
 259 parameters, recognizes when a certain output cannot be well predicted and then fixes an input  
 260 which is necessary to achieve higher accuracy.

261 In our experiments, all parameter fields are normalized before training to equilibrate differences  
 262 in parameter magnitudes and variation. If the dataset contains time-invariate parameters, such as  
 263 spatially varying, but temporally constant orography descriptors, these fields are concatenated to

264 the inputs to serve as additional information for the models. In particular, orography is used in  
265 the case of the WB dataset to enable the networks to learn dependencies between the parameters  
266 and land-/sea-scape, and, thus, enhance their inferencing skills. The quality of the reconstruction  
267 of all parameters is measured by a suitable loss function, e.g.  $L_1$  loss, and the model weights are  
268 optimized using standard backpropagation.

269 We monitor both the training and validation loss to avoid overfitting.

270 A network’s loss curve indicates how difficult it is for the network to achieve an accurate  
271 reconstruction depending on the current input and output parameters. I.e., depending on which  
272 parameters are used, the reconstruction error decreases more or less quickly. While saturation of  
273 both losses typically happens after 70 epochs, our experiments show that already after few epochs  
274 of training the reconstruction error clearly reveals the differences between different parameter  
275 combinations. In particular, when comparing the loss curves in these early stages with the loss  
276 curves after convergence, the relative behaviour of the networks does not change. This indicates  
277 that network training does not need to be performed until convergence, but can be stopped after  
278 few epochs to obtain an indication of the reconstruction quality. In particular, we consider loss  
279 values after five epochs of training, resulting in roughly one hour (for training 20 networks) on a  
280 low-size deep learning cluster with six mid-size GPUs to determine three input and three output  
281 parameters for the WB dataset, comprising six parameters. For the CSEns dataset, comprising nine  
282 different parameters, the proposed procedure requires roughly six hours for training 87 networks  
283 to determine the four input parameters that best predict the remaining five output parameters.

284 Compared to the V2V approach by Han et al. (2021b), the proposed strategy is computationally  
285 more expensive, yet it exhibits a number of advantages: Firstly, we obtain a more accurate measure  
286 of transferability, since our models are directly trained to reconstruct parameters. Second, the  
287 proposed approach is not constrained to selecting pairs of parameters, but can uncover multi-  
288 parameter relationships. Lastly, the method, in principle, enables to set a loss threshold for  
289 triggering the stopping of iterations. Due to normalization of the target parameters, this threshold  
290 can be interpreted as a measure of acceptable relative error, and is thus more accessible than the  
291 distance threshold in the latent-space of the original V2V algorithm.

292 **5. Datasets**

293 To validate our approach, we use the WeatherBench benchmark dataset, which has been proposed  
294 as a benchmark dataset for data-driven, medium-range climate prediction problems (Rasp et al.  
295 2020), as well as a convective-scale forecast ensemble, generated by Necker et al. (2020).

296 In both cases, the proposed neural network models receive as input an array of shape  $m \times H \times W$ ,  
297 with  $m$  the number of physical input parameters, and  $H \in \mathbb{N}$  and  $W \in \mathbb{N}$  denoting the spatial  
298 dimensions of the field. Assuming an initial number of  $n \in \mathbb{N}$  physical parameters, the model  
299 output is a field of shape  $(n - m) \times H \times W$ , which contains reconstructions of the parameters that  
300 have not been considered in the input (see Fig. 2).

301 *a. WeatherBench*

302 WeatherBench (WB) is based on ERA5 atmospheric reanalysis data (Hersbach et al. 2020),  
303 which are generated regularly at the European Center for Medium-Range Weather Forecasting  
304 (ECMWF) through data assimilation procedures, combining spatio-temporal numerical simulations  
305 and observation data.

306 To facilitate the accessibility to machine learning workflows and accelerate studies in weather  
307 prediction, WB provides regridded ERA5 reanalysis data on regular latitude-longitude grids with  
308 three different resolutions and 13 different pressure levels. The data is available hourly for 40 years  
309 from 1979 to 2018. For efficiency reasons, we consider a selection of 2D single-level fields with a  
310 resolution of  $1.40525^\circ$  in latitude and longitude, resulting in a domain size of  $128 \times 256$  vertices for  
311 global data. The selected physical parameters are 2 m-temperature (t2m), total cloud cover (tcc), u-  
312 and v-component of 10 m-wind (u10, v10), total precipitation (tp), and top-of-atmosphere incident  
313 solar radiation (tisr). For the WB data, we add orography height as an additional temporally steady  
314 predictor. Orography information is available at low memory cost, due to steadiness in time, and  
315 was found to support prediction accuracy of convolutional neural networks in earlier studies (e.g.,  
316 Höhlein et al. 2020). We use the 23 first years of WB during training. Of these, 20 years serve  
317 as training data for fitting the models, and three years are reserved for validation. The remaining  
318 years are left out for testing and visualization.

319 *b. Convective-scale ensemble*

320 The convective-scale ensemble simulation (CSEns), generated by (Necker et al. 2020), contains  
321 1000 runs of a 3D atmospheric dynamics model over a rectangular domain in central Europe. Data  
322 are stored on a regular grid with  $352 \times 250$  nodes, which corresponds to a horizontal grid spacing of  
323 3 km and allows the resolution of convective effects in the model dynamics. The simulation covers  
324 a time interval of six hours, with a period of one hour between successive time steps, and comprises  
325 30 levels in height. In the lower levels, some of the data are invalid due to grid cells falling below  
326 the level of the surface topography. Levels with missing values are omitted. 3D data are available  
327 for a total of 9 different parameters, which are temperature ( $tk$ ), u-, v-, and w- component of winds  
328 ( $u$ ,  $v$ ,  $w$ ), geopotential height ( $z$ ), relative humidity ( $rh$ ), mixing ratio of all hydro meteors ( $qh$ ),  
329 water vapor mixing ratio ( $qv$ ), and radar reflectivity ( $dbz$ ). Structural differences are observed  
330 not only between different parameters, but also between different timesteps and height levels of  
331 the same parameter. Specifically, the variation in the fields decreases with increasing distance  
332 from the earth surface, due to decreasing influence of boundary layer effects, and complexity  
333 increases with increasing simulation time due to a strengthening of convective activity. To enable  
334 a fair comparison between the CSEns and WB, we consider data only for the three lowermost  
335 levels without missing values, as well as the three latest time steps, which show the highest field  
336 complexity. We further split time-variate 3D fields both in time and height to obtain a sequence of  
337 plain 2D fields. We then consider data for 200 members for training, five members for validation,  
338 and the remaining members for testing and visualization.

339 **6. Network architecture**

340 In this study, we propose to train a deep convolutional neural network (CNN) architecture to  
341 predict a certain number of output parameters from a given set of input parameters. In general,  
342 deeper networks can have higher prediction quality, yet they can easily lead to convergence problems  
343 in the optimization process due to vanishing gradients (Glorot and Bengio 2010). In early layers of  
344 the network, gradient estimation causes an exponential decay of the gradient magnitudes, so that  
345 the parameters cannot change significantly in the training process. An efficient way to overcome  
346 this problem is to utilize short-cut or residual connections as in ResNet architectures (He et al.

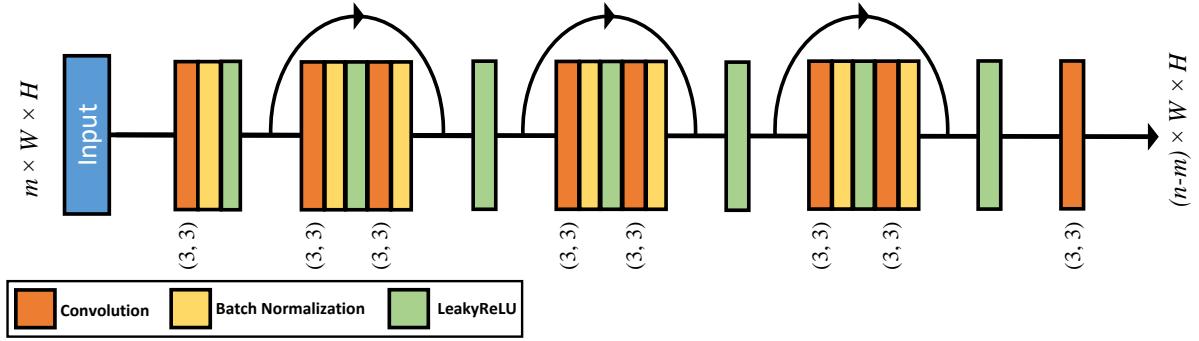


FIG. 3. Schematic of the used ResNet architecture. It consists of three residual blocks, each followed by a LeakyReLU activation function.  $m (< n)$  input fields are fed into the input convolution layer with 64 channels and kernel size of  $(3, 3)$ , a batch normalization layer, and LeakyReLU activation function. The network comprises three residual blocks and a final single convolution layer. The network predicts  $n - m$  output fields.

2016). In such architectures, outputs of earlier layers are added to the output of later layers, thus circumventing the accumulation of intermediate gradients.

In this study, we select a ResNet architecture with three residual blocks. A schematic representation is shown in Fig. 3. The input block of the model consist of a single convolution layer with kernel size of  $(3, 3)$ , 64 channels, batch normalization, and leaky rectified linear unit (LeakyReLU). We use input padding before each convolution layer. To account for periodic boundary conditions in the longitude direction of the WB dataset, we employ a periodic padding scheme in this dimension, and replication padding elsewhere. After that, there are three residual blocks and each block has two convolution layers with 64 channels and kernel size  $(3, 3)$ . Since the number of parameters grows with the kernel size, it is cost efficient to select a kernel of size 3. After the first convolution layer in the residual block, the network utilizes a batch normalization layer, a LeakyReLU layer, the second convolution layer, and another batch normalization layer. Batch normalization is used to achieve improved stability and convergence (Ioffe and Szegedy 2015). After each residual block, a LeakyReLU activation function guarantees non-linearity of the mapping. The final layer is a single convolution layer with kernel size of  $(3, 3)$  and  $n - m$  output channels.

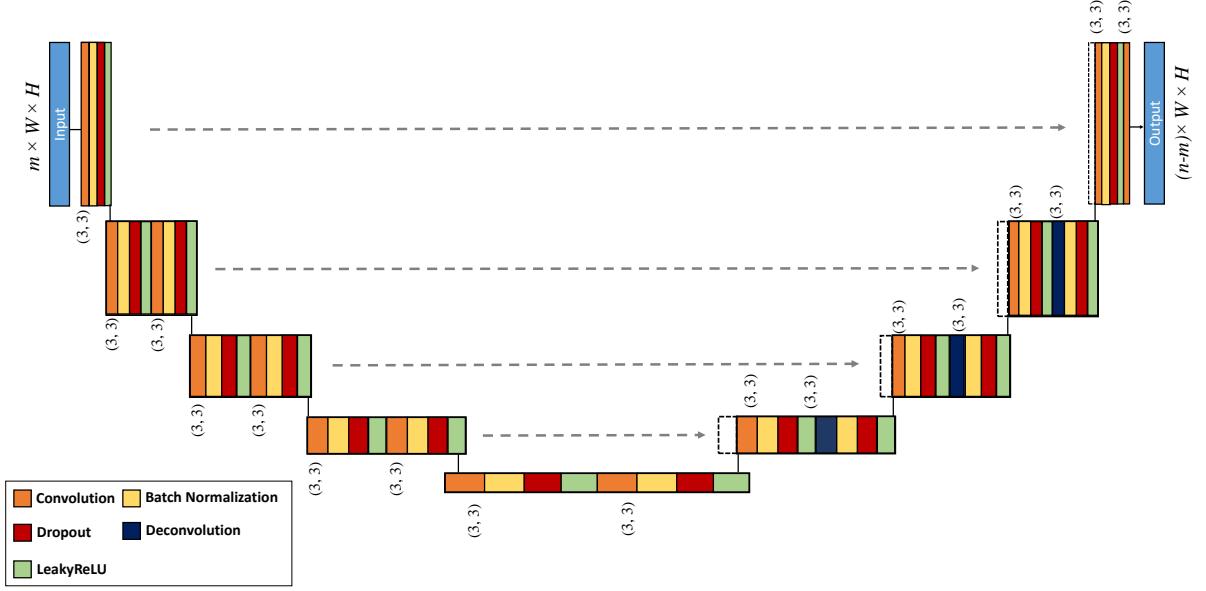
As an alternative to the ResNet architecture, we also analysed the potential of a UNet architecture (Ronneberger et al. 2015) for loss-based parameter selection.

368 In contrast to the ResNet architecture, which operates on a single spatial scale throughout the  
369 whole architecture, the UNet architecture allows for the extraction of features on multiple spatial  
370 scales, which offers the possibility to learning a wider range parameter relationships. The UNet  
371 consists of two symmetric branches, which give it the characteristic u-shape, as seen in Fig. 4. In  
372 the encoding branch, the data are encoded into an abstract reduced feature representation, and in the  
373 decoding path the feature representations are then decoded to reconstruct the predicted fields at the  
374 target resolution. During the encoding step, the resolution is iteratively reduced, and the number  
375 of feature channels is increased at the same time. In the decoding step, while reducing the number  
376 of feature channels, the features are super-sampled to a higher resolution. The paths are connected  
377 by skip connections, which concatenate feature channels from the encoder with corresponding  
378 features from the decoder, in order to precisely preserve and localize the information in the data  
379 that could be lost in the encoding stage. The most bottom layer of the UNet, i.e., the bottleneck  
380 layer, enforces the model to learn a compact representation of the input containing the globally  
381 most relevant information to recover it.

382 In our experiments, the UNet architecture did not improve the reconstruction quality significantly,  
383 yet increased the training time due to its higher computational complexity. A sample of the  
384 reconstruction quality of the UNet architecture is shown in Fig. A2. Nevertheless, we found that  
385 ResNet and UNet seem to learn different mappings internally, which we discuss in more detail in  
386 section 7 c.

387 The presented architectures have been designed through empirical experimentation, trading off  
388 model flexibility and reconstruction quality against applicability to diverse datasets and computa-  
389 tional efficiency. Especially for data fields on spherical geometries, more sophisticated network  
390 designs exist, which account for domain periodicity, grid distortion and rotation equivariance more  
391 accurately, see, e.g., the survey by Cao et al. (2020). Such architectures, however, come at higher  
392 computational complexity or require careful data-specific selection of hyper-parameters to achieve  
393 better performance than standard CNNs, and are thus not considered in the present study.

401 The error between target fields and predictions is measured in terms of  $L_1$  distance, which we  
402 prefer over  $L_2$  due to empirically less pronounced suppression of outlying predictions.



394 FIG. 4. Schematic of the used UNet architecture. The contracting branch is comprised of three convolutional  
 395 blocks, each consisting of two convolution layers with subsequent batch normalization, dropout, and a  
 396 LeakyReLU activation function. The expansive branch includes three deconvolution blocks, each consisting of  
 397 one convolution, and one deconvolution layer. Each of these layers is followed by a batch normalization layer,  
 398 dropout, and LeakyReLU activation function.  $m$  ( $< n$ ) input fields are fed into the input block, which contains  
 399 a single convolution layer with 64 channels and kernel size of (3, 3), followed by batch normalization layer,  
 400 dropout, and LeakyReLU activation function. The number of reconstructed fields is  $n - m$ .

## 403 7. Experiments

404 In an exhaustive ablation, we demonstrate the feasibility and reliability of our approach using the  
 405 WB reanalysis dataset as a use case. The results of applying the proposed strategy to the CSEns  
 406 dataset are shown in the appendix.

407 Via this ablation study, we aim to answer the following questions:

- 408 1. Which is the minimal set of input parameters from which the remaining parameters can be  
 409 reconstructed accurately? This number indicates how aggressively the initial parameter set  
 410 can be reduced.
- 411 2. Over which geographic regions do parameters strongly affect the network's prediction quality?

412 For answering the first question, we use loss-based parameter selection via the ResNet architec-  
413 ture. Both architectures are used subsequently to answer the second question by visualizing the  
414 sensitivity of the local prediction accuracy to regional changes of the input parameters.

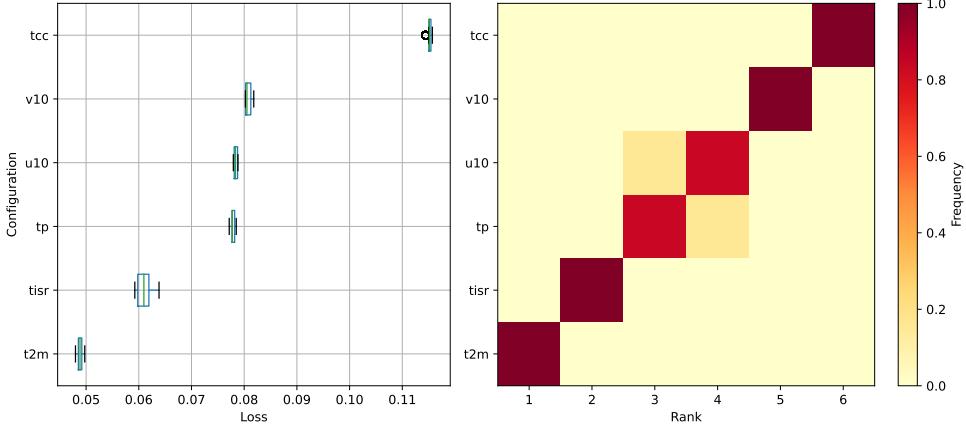
415 *a. Validation of the extended V2V approach*

416 To validate the reliability and reproducibility of the proposed loss-based parameter selection,  
417 we train networks for different parameter configurations multiple times with different random  
418 weight initializations, and compare the order of the observed losses after five training epochs.  
419 For brevity, we first show results only for the case  $m = 1$ , which results in six different parameter  
420 configurations. For every configuration, we train 10 models and sample model ensembles by  
421 randomly picking one of the 10 models for each configuration. For each sample, we then rank the  
422 model configurations according to the observed loss value after five training epochs, and assess  
423 the consistency of the ranking order among different samples. Fig. 5 illustrates the observed loss  
424 statistics. We find that the separation in loss magnitude between different parameter configurations  
425 is typically larger than the variance of losses for each configuration (see Fig. 5, left). As a result,  
426 the ranking of losses is consistent between different runs. This is seen in the heat chart in Fig. 5  
427 (right), which visualizes the frequency of how often a particular loss rank is observed for each of  
428 the parameters, and suggests an almost perfect one-to-one mapping between parameters and ranks.  
429 Both charts together confirm that our loss-tracking approach constitutes a reproducible criterion for  
430 selecting parameter configurations. Nevertheless, we observe that clustering of losses may occur,  
431 i.e. different configurations may result in very similar loss statistics (e.g., parameters tp, u10 and  
432 v10 in Fig. 5, left).

433 In such cases, the decisions of our algorithm may differ between training runs. Due to the overall  
434 small variation in losses per configuration, we conjecture that all of the possible outcomes are  
435 equally well-suited for further evaluations.

439 *b. Ablation study*

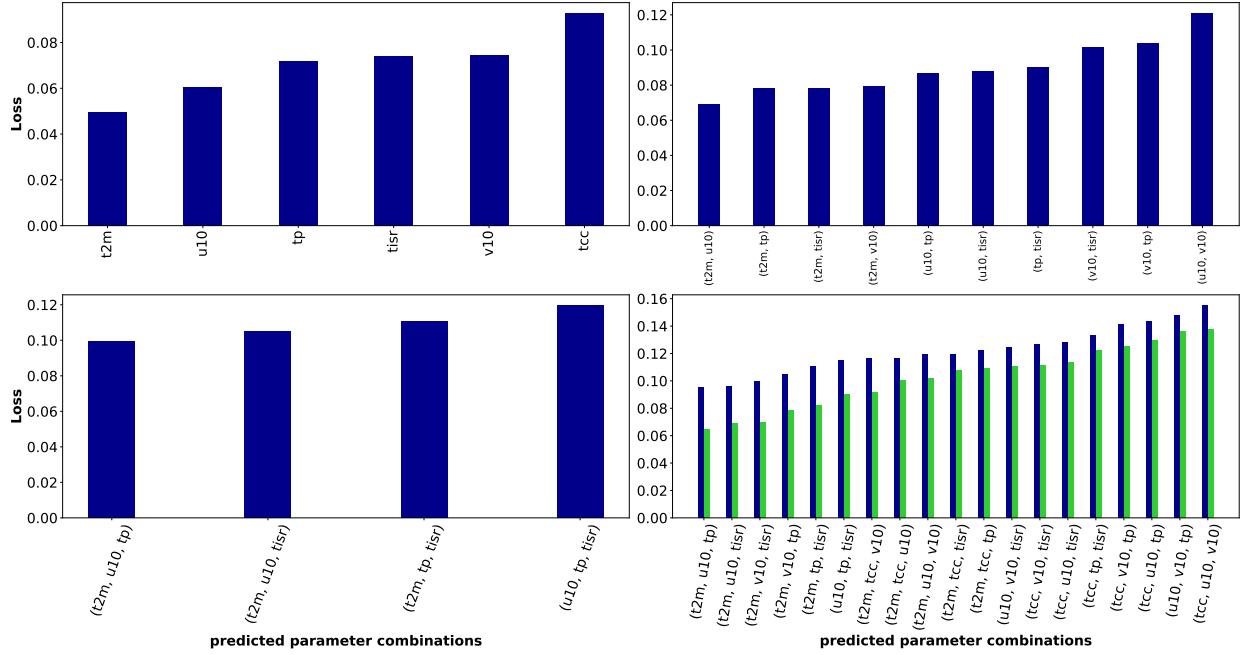
440 For the WB data comprising six parameters, we start with applying the loss-based selection  
441 procedure to predict one masked out parameter using the remaining five input parameters. This  
442 number is then increased to two and finally three predicted parameters, with four and three input



436 FIG. 5. Distribution and ranking of losses for different network configurations. Networks are trained for  
 437 different parameter configuration ( $m = 1$ , i.e. one parameter is left out and is to be predicted) with different  
 438 weight initializations. Left: Box plot of the loss statistics. Right: Heat map of the observed ranking order.

443 parameters, respectively. This means that during the first iteration six networks, then  $\binom{5}{2}$  networks  
 444 and finally  $\binom{4}{3}$  networks are trained. We do not go beyond three masked out parameters, since  
 445 significantly reduced reconstruction quality is observed in this case.

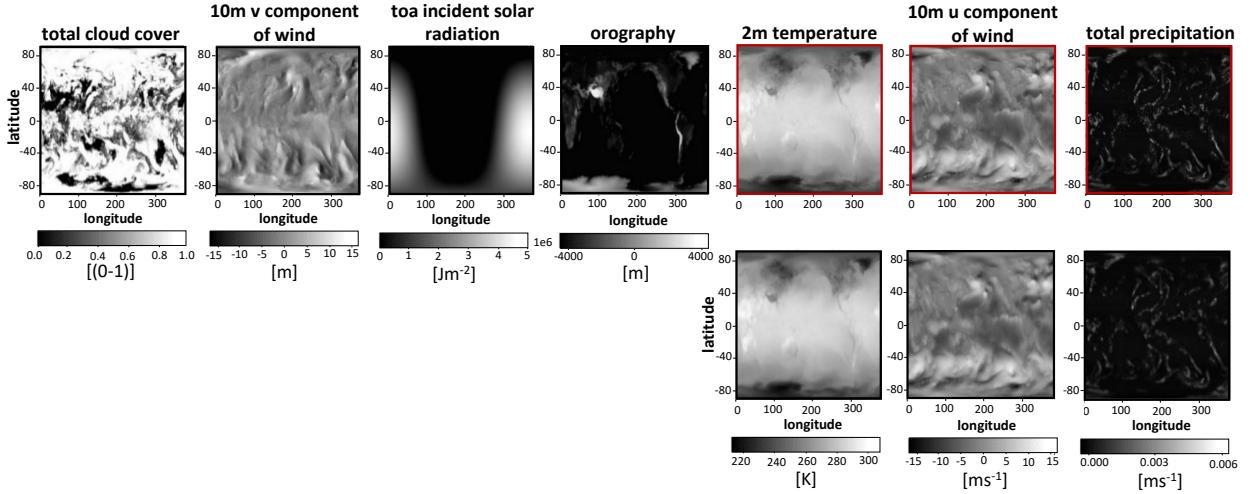
450 To justify our decision to use the network losses after five epochs as indicators of the difficulty  
 451 to predict a certain parameter or parameter combination, we analyze the loss values for five and  
 452 70 epochs of training of all networks that were trained. Fig. 6 shows the losses of all networks  
 453 trained for five epochs by the proposed loss-based selection approach (top and bottom left charts),  
 454 and the losses of all  $\binom{6}{5}$  possible networks trained for five epochs (dark blue bars in bottom right).  
 455 The loss values indicate that the loss-based selection approach finds the parameter combination  
 456 yielding the lowest loss. Note that this is also confirmed for the CSEns dataset, as shown in Figs. B2  
 457 and B3 in the appendix. As shown by the overlayed loss values of the networks trained for 70  
 458 epochs (green bars in bottom right), training for five epochs shows very similar relative differences  
 459 between different parameter combinations. Also this result is confirmed by the comparison of the  
 460 loss values for five and 70 of the CSEns dataset. When only one parameter is predicted from the  
 461 remaining five parameters (plus orography), it can be seen that 2 m-temperature and total cloud  
 462 cover, respectively, are the parameters that are easiest and most difficult to reconstruct. Thus, total  
 463 cloud cover is the first parameter that is fixed in the input.



446 Fig. 6. Bar charts showing the losses of all networks trained for 3-to-3 parameter transfer with the WB dataset  
447 using the proposed iterative loss-based approach (top left: first iteration, top right: second iteration, bottom left:  
448 third iteration), and of all possible  $\binom{6}{3}$  networks (bottom right). Blue bars represent losses after five epochs of  
449 training, green bars indicate losses after 70 epochs.

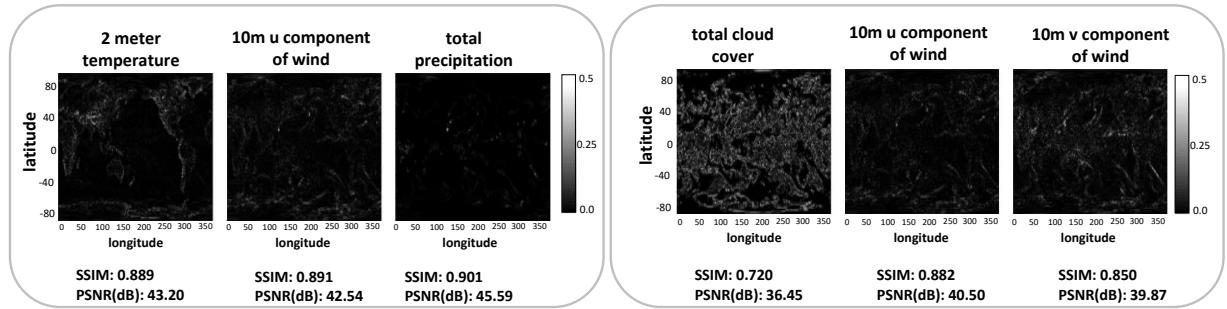
464 Fig. 7 shows the initial parameter fields (including orography) and the reconstruction results  
465 of the three parameters that have been masked out by the loss-based procedure. For comparison,  
466 the reconstruction results of the three worst parameter combinations are shown in A1 in the  
467 appendix. Notably, when all  $\binom{6}{3}$  parameter combinations are evaluated, the very same combination  
468 is determined.

475 In Fig 8, for the selected parameter combination the resulting pixel-wise differences between the  
476 reconstructions and the initial parameter fields are shown. The quality of the reconstructed fields  
477 is measured using the image statistics SSIM (Wang et al. 2004) and the peak signal to noise ratio  
478 (PSNR), with the initial parameter fields as references. It can be seen that even when one half of  
479 the parameters are masked out, they can still be reconstructed at high accuracy by the network.  
480 In addition, also the reconstruction quality that is achieved by the worst parameter combination  
481 is shown, i.e., the parameter combination yielding the highest loss of all possible  $\binom{6}{3}$  parameter  
482 combinations. The results indicate the importance of a suitable procedure for finding the best



469 FIG. 7. Reconstruction results for the WB dataset when the network is trained to predict three parameter fields  
 470 from three input fields and orography. Top: The initial parameter fields. A red outline indicates those fields the  
 471 network has learned to predict from the others. Bottom: Predicted parameter fields.

483 parameter combination. The pixel-wise error plots indicate significantly different reconstruction  
 quality between the best and worst parameter set.



472 FIG. 8. Pixel-wise differences between the the initial and predicted fields when using the best (left) and worst  
 473 (right) parameter combination. Per-pixel values are scaled by a factor of 10 for better visibility. Corresponding  
 474 SSIM and PSNR (dB) values are given below each image.

484

### 485 c. Feature analysis

486 While the potential of the selected network architecture for V2V can be concluded from the  
 487 results of the ablation study, no information can be drawn about what kind of dependencies are

488 exploited by the networks. To shed light on this aspect, we use layer-wise relevance propagation  
 489 (LRP) to localize the sensitivity of the reconstruction results to changes in the input parameter  
 490 fields.

491 In its original form, LRP has been introduced as an explainability algorithm for image clas-  
 492 sification models (Bach et al. 2015), which is achieved by combining neuron activations and  
 493 back-propagated gradient information to highlight image regions that exert a strong effect on the  
 494 classifier output. LRP, thereby, builds on the concept of pixel-wise decomposition of the classifier  
 495 score. I.e., given a classification score mapping of the form  $f : \Omega \rightarrow \mathbb{R}$ , where  $\Omega \subseteq \mathbb{R}^{m \times H \times W}$  is the  
 496 input domain (e.g., the space of images with  $H \times W$  pixels and  $m$  channels per pixel), and  $f(x) > 0$   
 497 ( $< 0$ ) indicates evidence for presence (absence) of a particular feature, LRP attempts to find a set  
 498 of relevance values  $R_{kij} \in \mathbb{R}$  associated with the pixel values, such that the classification score can  
 499 be approximated as

$$f(x) \approx \sum_{k=0}^{m-1} \sum_{i=0}^{H-1} \sum_{j=0}^{W-1} R_{kij}, \quad (1)$$

500 and  $R_{kij} > 0$  ( $< 0$ ) indicates that pixel channel  $k$  at position  $(i, j)$  contributes evidence in favor  
 501 of (against) the presence of the feature in question. In the case of deep neural networks, which  
 502 are composed of linear transformations with element-wise activation functions, suitable relevance  
 503 values can be computed via iterative relevance back-propagation, subject to propagation rules (Bach  
 504 et al. 2015).

505 Deviating from the setting of standard LRP, the input of V2V models is not an image, but  
 506 a multi-dimensional array, representing a multi-parameter field, and the output is not a uni-  
 507 variate classification score, but a multi-dimensional multi-parameter field. The difference in input  
 508 modalities is only of limited importance, since the multi-parameter fields can be interpreted directly  
 509 as multi-channel images. However, the complexity of the model output prevents straight forward  
 510 application of standard LRP. We therefore propose to use an adapted variant of LRP to gain insight  
 511 into spatio-temporal relevance and correlation patterns between model predictions and inputs.

512 Given a model mapping of the form  $f : \Omega \rightarrow \mathbb{R}^{(n-m) \times H \times W}$ , we propose adding an additional  
 513 selector layer  $s : \mathbb{R}^{(n-m) \times H \times W} \rightarrow \mathbb{R}$  at the end of the model, such the output of the combined model,  
 514  $s(f(x)) \in \mathbb{R}$ , admits an additive decomposition according to Eq. (1), and can thus be further  
 515 analyzed using standard LRP. Possible choices for  $s$  include summarization operators, such as

516 global (or local) averaging of field values or deviation measures, or selection operations, which  
 517 select single pixels and output channels for computing LRP relevances. Depending on the choice  
 518 of selector layer, different aspects of the input-output relationship can be investigated. For instance,  
 519 a selector function returning the mean value of the output channel  $0 \leq c < m$  inside a region defined  
 520 by the pixel set  $I \subseteq \{(i, j) : 0 \leq i < H, 0 \leq j < W\}$ , i.e.

$$s_I^{(c)}(x) := \langle [x]_{kij} \rangle_{(i,j) \in I, k=c}, \quad (2)$$

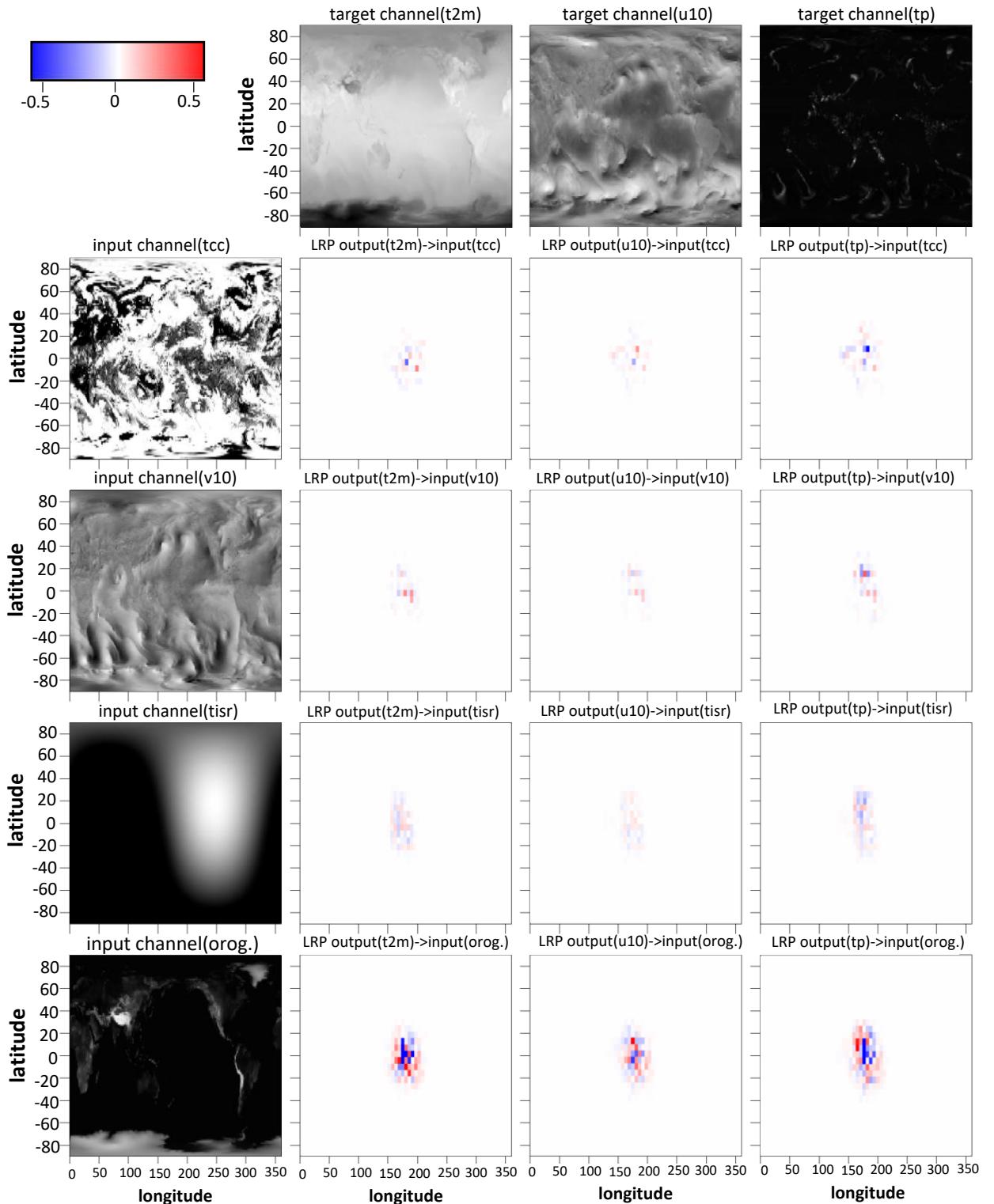
521 where  $[x]_{kij}$  denotes selection of the element at position  $(k, i, j)$  in the array  $x$ , yields positive  
 522 relevance for input regions. This causes an increase of the averaged quantity according to Eq. (1).  
 523 In contrast, functions of the form

$$\delta_I^{(c)}(x; x_0, p) := \langle |[x - x_0]_{kij}|^p \rangle_{(i,j) \in I, k=c}, \quad (3)$$

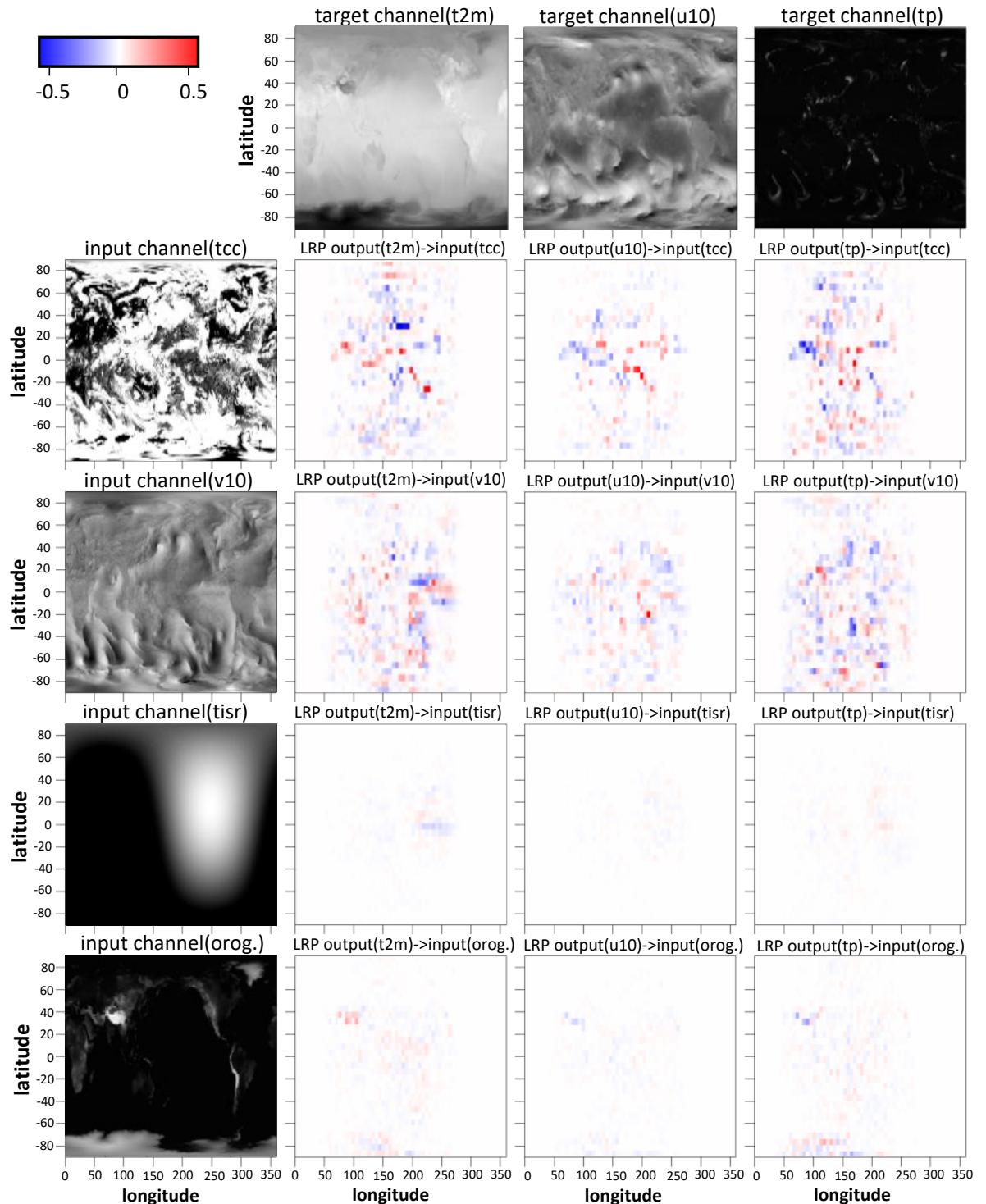
524  $p > 0$ , yield positive relevance for regions which increase the deviation between the model predic-  
 525 tion  $x$  and a certain reference prediction  $x_0$  within the region  $I$ .

526 Figs. 9 and 10 show relevance maps for the global atmospheric situation on May 15, 2004,  
 527 08h, as seen in the WB dataset, using the ResNet (Fig. 9) and the UNet model (Fig. 10). We  
 528 employ an absolute-difference-based selector function with a focus on single pixel deviations  
 529 of the predicted quantities from the respective target value, i.e.  $\delta_I^{(c)}(x; x_0, 1)$  with  $I = \{(i^*, j^*)\},$   
 530  $0 \leq i^* < H, 0 \leq j^* < W$ , and  $x_0$  denoting the target field. This allows drawing information about what  
 531 parts in the data push the separate prediction channels away from the actual target value. Figures  
 532 are shown for  $(i^*, j^*) = (63, 127)$ , which corresponds to the center pixel in the image, located at  
 533  $0^\circ\text{N } 180^\circ\text{E}$ . Relevance maps for other dates and pixel indices look similar. All output channels  
 534 are treated separately, yielding a matrix of relevance maps, which visualize relationships between  
 535 channel-wise prediction errors and model inputs. The computation of relevances is based on the  
 536 LRP implementation available in the Captum model interpretability library for Pytorch Kokhlikyan  
 537 et al. (2020).

538 While the relevance patterns appear noisy in both cases, the structure of the relevance maps  
 539 differs significantly between the model architectures, despite being trained on the same task and  
 540 with the same set of training data. The ResNet architecture favors relevance distributions which



526 FIG. 9. LRP relevance maps with deviation-based selector function for the ResNet model in the best input-  
 527 output configuration wrt. the proposed selection procedure. Timestamp of data sample: May 15, 2004, 08h.



528 FIG. 10. LRP relevance maps with deviation-based selector function for the UNet model in the best input-output  
529 configuration wrt. the proposed selection procedure. Timestamp of data sample: May 15, 2004, 08h

545 are concentrated around the reference location. This is consistent with the inductive bias of  
546 the architecture, which arises from the use of convolution layers with small kernel sizes (see  
547 section 6). Positive and negative relevances appear to be distributed randomly throughout the map,  
548 but significant differences are observed in the magnitude of the relevances. Relevance values wrt.  
549 orography possess larger magnitude in both positive and negative orientation than the remaining  
550 parameters. For the cloud-cover input, relevance values are concentrated on a small number of  
551 pixels, which obtain a relevance with notably higher amplitude than that of surrounding pixels.  
552 Similarly, the large degree of variability in the relevance maps makes it difficult to identify. For  
553 the UNet, relevance values are distributed over a larger spatial domain and relevance magnitude  
554 is largest for the cloud cover field and the v-component of the wind field. Likely, this is caused  
555 by the multi-scale properties of the UNet architecture, and confirms that the UNet manages to  
556 learn features on larger spatial scales. Notably, the distribution of relevance values also displays  
557 stronger spatial correlations, which might suggests that the model learns to pay attention to spatially  
558 coherent features in the data. Also, in contrast to ResNet, the relevance of the orography field  
559 is smaller. Intuitively, this relevance attribution appears more understandable, since the selected  
560 reference pixel corresponds to a location in the mid of the Pacific Ocean, where the impact of  
561 orography on physical processes should be weak.

562 A prominent feature of the WB dataset is the temporal coherence of subsequent samples, which  
563 is determined by the day-night cycle, as well as the seasonal cycle. To assess the stability of the  
564 relevance maps, as well as the impact of the day-night cycle on the model mapping, we show  
565 two additional relevance maps for the UNet model applied to data samples from May 15, 2004,  
566 12h and 20h in Figs. A3 and A4 in the appendix. The figures show that the relevance maps for  
567 08h and 20h look very similar. In particular, clusters of spatially coherent regions of positive or  
568 negative relevance are preserved, which suggests a stability and coherence in the visual structure  
569 of the relevance maps. The maps for 12h deviate slightly and show larger relevance as for the 2 m-  
570 temperature with respect to the top-of-atmosphere incoming solar radiation, which is consistent  
571 with physical intuition.

572 Overall, we conclude that the different architectures learn distinct mappings, despite being trained  
573 on the same task and achieving similar prediction accuracy. Yet, we find that some aspects of the  
574 dependency structure can be partly reverse engineered via thorough investigation of the relevance

maps. Similar statements apply to the relevance maps for models trained on the CSEns dataset. Exemplary relevance maps for this dataset are shown in Fig. B8. A more detailed analysis of the derived relevance maps, as well as the study of more specific meteorological events and weather situations at selected times and locations, is however beyond the scope of this paper and will be addressed in future work.

## 580 8. Conclusion

581 We have introduced an alternative way to perform deep learning-based variable-to-variable  
582 transfer. Instead of building upon the similarity of latent-space representations of parameter fields  
583 to determine transferable parameter pairs, we train a network using different transfer scenarios and  
584 select the best parameter setting. In this way, we give more flexibility to the network to exploit  
585 inter-parameter relationships, i.e., to learn parameter combinations for improved transfer. When  
586 applied to weather data, the results indicate that even for parameter fields with vastly different  
587 distributions, one half of the parameter set is sufficient to reconstruct the other half. This allows  
588 saving bandwidth in in-situ settings, and can help to more aggressively compress multi-parameter  
589 simulation data. As shown in Fig. B9 in the appendix, V2V transfer cannot compete with classical  
590 lossy data compression schemes in terms of compression rate, yet it may effectively support such  
591 schemes when the structure of the relevance maps generated via LRP is exploited to select spatially  
592 varying bitrates according to the importance of the data values.

593 We have further analyzed the regional parameter structures that have the most significant effect  
594 on the reconstruction quality. By using an extension to layer-wise relevance propagation (LRP),  
595 we were able to determine regions over which the field values have a large effect on the local  
596 reconstruction accuracy. LRP results demonstrate that different model architecture learn different  
597 mapping functions, depending on the inductive bias of the used architecture. In our study, the use  
598 of the UNet model led to more physically interpretable relevance maps, while the mappings learned  
599 by the ResNet architecture are constrained in learning spatial dependencies due to the construction  
600 of the network architecture. Information like this may help in operational model applications to  
601 gain a better understanding of model-driven inference procedures and increase trustworthiness of  
602 data-driven model predictions.

603 In the future, we will shed light on the use of the proposed V2V approach with 3D and especially  
604 large forecast ensembles. In our current use cases, all parameter fields show rather low mutual  
605 similarities, and, thus, one can expect our approach to perform even more effective once parameter  
606 fields with certain similarities and more pronounced spatial relationships are given, like ensemble  
607 simulations. One specific task we envision is to analyse the representativeness of the single  
608 members captured by a Grand Ensemble, by using V2V to reconstruct an as small as possible  
609 subset of all members capturing the full ensemble spread. This can facilitate guidance towards  
610 weather situations that are under- or over-represented in the ensemble, and reveal situations which  
611 are intrinsically difficult to resolve. Furthermore, we intend to consider the temporal evolution  
612 of the fields to improve the reconstruction at a certain time, i.e., by letting the network train on  
613 multiple timesteps from the past.

614 Finally, together with meteorologists and climatologists we intend to further analyse the sensitiv-  
615 ity maps that have been derived via LRP. Such an analysis includes the extraction of specific local  
616 weather events such as jet-cores or fronts, and to set them into relation to the regions that have been  
617 deemed important for achieving high reconstruction accuracy. A limitation of the current LRP  
618 approach lies in the necessity of selecting reference locations, for which "point-to-field" relevance  
619 maps shall be computed. In exploratory data analysis tasks, it might be non-trivial to make sen-  
620 sible decisions about which locations to look at in a first place. We therefore plan on refining the  
621 LRP-based analysis procedures to detect regions of high impact in an automated fashion and with  
622 a more global view to enable the interactive exploration "field-to-field" relevance relations. In a  
623 similar line of reasoning, we intend to include the time dimension in the analysis, e.g., by using  
624 temporal coherence and recurrence in the data to reduce the noise level of the derived LRP maps  
625 via temporal filtering or climatological summarization of relevances. Further efforts will be put on  
626 the investigation of alternative mechanisms for pursuing a sensitivity analysis, focusing more on  
627 spatial as well as temporal relationships between different parameters.

628 *Acknowledgments.* This study has been done within the subproject “Visualization of coherence  
629 and variation in meteorological dynamics“ of the Transregional Collaborative Research Center  
630 SFB/TRR 165 “Waves to Weather“ funded by the German Research Foundation (DFG).

631 *Data availability statement.* WeatherBench dataset (Rasp et al. 2020) is publicly available at  
632 <https://github.com/pangeo-data/WeatherBench>. Access to the convective-scale ensemble  
633 data can be requested from the authors of the dataset, Necker et al. (2020). The code for the  
634 experiments is made publically available at [https://github.com/FatemehFarokhmanesh/  
635 DNN-based-Parameter-Transfer-in-Meteorological-Data.git](https://github.com/FatemehFarokhmanesh/DNN-based-Parameter-Transfer-in-Meteorological-Data.git).

## 636 References

- 637 Bach, S., A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek, 2015: On pixel-wise  
638 explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS one*,  
639 **10** (7), e0130140.
- 640 Cao, W., Z. Yan, Z. He, and Z. He, 2020: A comprehensive survey on geometric deep learning.  
641 *IEEE Access*, **8**, 35 929–35 949.
- 642 Cheng, J., Q. Kuang, C. Shen, J. Liu, X. Tan, and W. Liu, 2020: Reslap: Generating high-resolution  
643 climate prediction through image super-resolution. *IEEE Access*, **8**, 39 623–39 634.
- 644 Glorot, X., and Y. Bengio, 2010: Understanding the difficulty of training deep feedforward neural  
645 networks. *Proceedings of the thirteenth international conference on artificial intelligence and*  
646 *statistics*, JMLR Workshop and Conference Proceedings, 249–256.
- 647 Guo, L., S. Ye, J. Han, H. Zheng, H. Gao, D. Z. Chen, J.-X. Wang, and C. Wang, 2020: SSR-  
648 VFD: Spatial super-resolution for vector field data analysis and visualization. *2020 IEEE Pacific*  
649 *Visualization Symposium (PacificVis)*, IEEE Computer Society, 71–80.
- 650 Han, J., and C. Wang, 2019: TSR-TVD: Temporal super-resolution for time-varying data analysis  
651 and visualization. *IEEE transactions on visualization and computer graphics*, **26** (1), 205–215.
- 652 Han, J., and C. Wang, 2020: SSR-TVD: Spatial super-resolution for time-varying data analysis and  
653 visualization. *IEEE Transactions on Visualization and Computer Graphics*.

- 654 Han, J., H. Zheng, D. Z. Chen, and C. Wang, 2021a: STNet: An end-to-end generative framework  
655 for synthesizing spatiotemporal super-resolution volumes. *IEEE Transactions on Visualization*  
656 and Computer Graphics
- 657 Han, J., H. Zheng, Y. Xing, D. Z. Chen, and C. Wang, 2021b: V2v: A deep learning approach to  
658 variable-to-variable selection and translation for multivariate time-varying data. *IEEE Transac-*  
659 *tions on Visualization and Computer Graphics*, **27**, 1290–1300.
- 660 He, K., X. Zhang, S. Ren, and J. Sun, 2016: Deep residual learning for image recognition.  
661 *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- 662 Helbig, C., L. Bilke, H.-S. Bauer, M. Böttinger, and O. Kolditz, 2015: Meva-an interactive  
663 visualization application for validation of multifaceted meteorological data with multiple 3d  
664 devices. *PloS one*, **10** (4), e0123811.
- 665 Hersbach, H., and Coauthors, 2020: The era5 global reanalysis. *Quarterly Journal of the Royal*  
666 *Meteorological Society*, **146** (730), 1999–2049.
- 667 Höhlein, K., M. Kern, T. Hewson, and R. Westermann, 2020: A comparative study of convolu-  
668 tional neural network models for wind field downscaling. *Meteorological Applications*, **27** (6),  
669 <https://doi.org/https://doi.org/10.1002/met.1961>.
- 670 Ioffe, S., and C. Szegedy, 2015: Batch normalization: Accelerating deep network training by  
671 reducing internal covariate shift. *International conference on machine learning*, PMLR, 448–  
672 456.
- 673 Kokhlikyan, N., and Coauthors, 2020: Captum: A unified and generic model interpretability  
674 library for pytorch. *arXiv preprint arXiv:2009.07896*.
- 675 Necker, T., S. Geiss, M. Weissmann, J. Ruiz, T. Miyoshi, and G.-Y. Lien, 2020: A convective-scale  
676 1,000-member ensemble simulation and potential applications. *Quarterly Journal of the Royal*  
677 *Meteorological Society*, **146** (728), 1423–1442.
- 678 Pouliot, D., R. Latifovic, J. Pasher, and J. Duffe, 2018: Landsat super-resolution enhancement  
679 using convolution neural networks and sentinel-2 for training. *Remote Sensing*, **10** (3), 394.

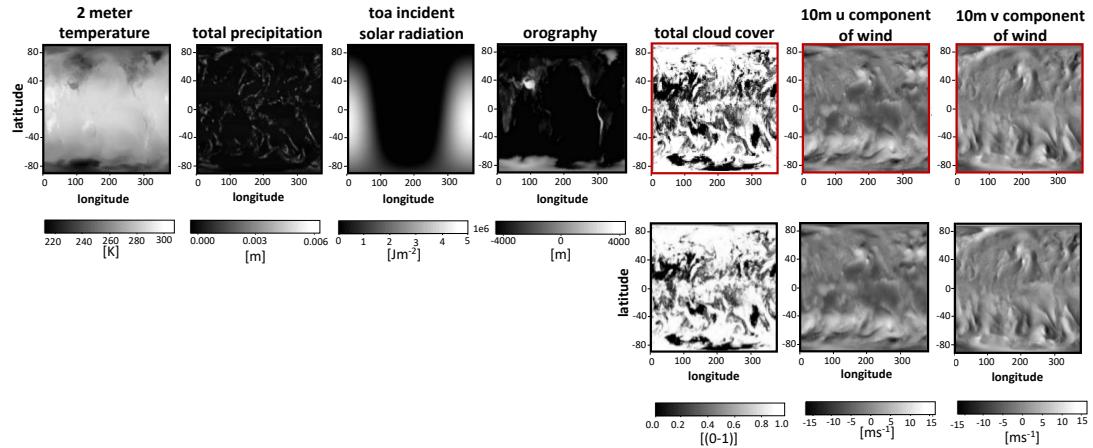
- 680 Rasp, S., P. D. Dueben, S. Scher, J. A. Weyn, S. Mouatadid, and N. Thuerey, 2020: Weatherbench:  
681 a benchmark data set for data-driven weather forecasting. *Journal of Advances in Modeling*  
682 *Earth Systems*, **12** (11), e2020MS002 203.
- 683 Reichstein, M., G. Camps-Valls, B. Stevens, M. Jung, J. Denzler, N. Carvalhais, and Coauthors,  
684 2019: Deep learning and process understanding for data-driven earth system science. *Nature*,  
685 **566** (7743), 195–204.
- 686 Röber, N., and J. F. Engels, 2019: In-situ processing in climate science. *International Conference*  
687 *on High Performance Computing*, Springer, 612–622.
- 688 Rodrigues, E. R., I. Oliveira, R. Cunha, and M. Netto, 2018: Deepdownscale: a deep learning  
689 strategy for high-resolution weather forecast. *2018 IEEE 14th International Conference on e-*  
690 *Science (e-Science)*, IEEE, 415–422.
- 691 Ronneberger, O., P. Fischer, and T. Brox, 2015: U-net: Convolutional networks for biomedical  
692 image segmentation. *International Conference on Medical image computing and computer-*  
693 *assisted intervention*, Springer, 234–241.
- 694 Sato, T., O. Tatebe, and H. Kusaka, 2019: In-situ data analysis system for high resolution meteoro-  
695 logical large eddy simulation model. *Proceedings of the 6th IEEE/ACM International Conference*  
696 *on Big Data Computing, Applications and Technologies*, 155–158.
- 697 Serifi, A., T. Günther, and N. Ban, 2021: Spatio-temporal downscaling of climate data using  
698 convolutional and error-predicting neural networks. *Frontiers in Climate*, **3**, <https://doi.org/10.3389/fclim.2021.656479>.
- 700 Toderici, G., D. Vincent, N. Johnston, S. Jin Hwang, D. Minnen, J. Shor, and M. Covell, 2017:  
701 Full resolution image compression with recurrent neural networks. *Proceedings of the IEEE*  
702 *conference on Computer Vision and Pattern Recognition*, 5306–5314.
- 703 Treib, M., F. Reichl, S. Auer, and R. Westermann, 2012: Interactive editing of gi-  
704 gasample terrain fields. *Computer Graphics Forum (Proc. Eurographics)*, **31** (2), 383–  
705 392, <https://doi.org/10.1111/j.1467-8659.2012.03017.x>, URL <http://diglib.eg.org/EG/CGF/volume31/issue2/v31i2pp383-392.pdf>.

- 707 van der Maaten, L., and G. Hinton, 2008: Visualizing data using t-SNE. *Journal of Machine*  
708 *Learning Research*, **9**, 2579–2605, URL <http://www.jmlr.org/papers/v9/vandermaaten08a.html>.
- 709 Wang, C., H. Yu, R. W. Grout, K.-L. Ma, and J. H. Chen, 2011: Analyzing information trans-  
710 fer in time-varying multivariate data. *2011 IEEE Pacific Visualization Symposium*, 99–106,  
711 <https://doi.org/10.1109/PACIFICVIS.2011.5742378>.
- 712 Wang, Z., A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, 2004: Image quality assessment: from  
713 error visibility to structural similarity. *IEEE transactions on image processing*, **13** (4), 600–612.
- 714 Weiss, S., M. Chu, N. Thuerey, and R. Westermann, 2019: Volumetric isosurface rendering  
715 with deep learning-based super-resolution. *IEEE Transactions on Visualization and Computer*  
716 *Graphics*, 1–1.
- 717 Weiss, S., M. Işık, J. Thies, and R. Westermann, 2020: Learning adaptive sampling and recon-  
718 struction for volume visualization. *IEEE Transactions on Visualization and Computer Graphics*,  
719 1–1.
- 720 Zhou, Z., Y. Hou, Q. Wang, G. Chen, J. Lu, Y. Tao, and H. Lin, 2017: Volume upscaling with  
721 convolutional neural networks. *Proceedings of the Computer Graphics International Conference*,  
722 1–6.

## APPENDIX A

## Supplementary Visualizations for the WeatherBench Dataset

725 The appendix provides supplementary figures illustrating specific aspects of V2V transfer in the  
 726 first dataset, the WeatherBench reanalysis (WB) dataset.



727 FIG. A1. Reconstruction results for the WeatherBench dataset when the network is trained to predict three  
 728 parameter fields (worst combination) from four input fields. Top: The initial parameter fields. A red outline  
 729 indicates those fields the network has learned to predict from the others. Bottom: Predicted parameter fields.

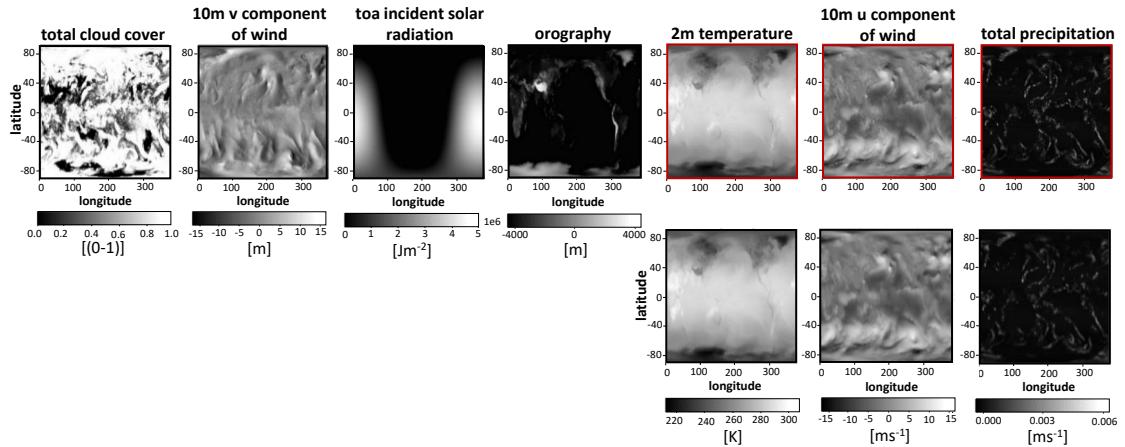
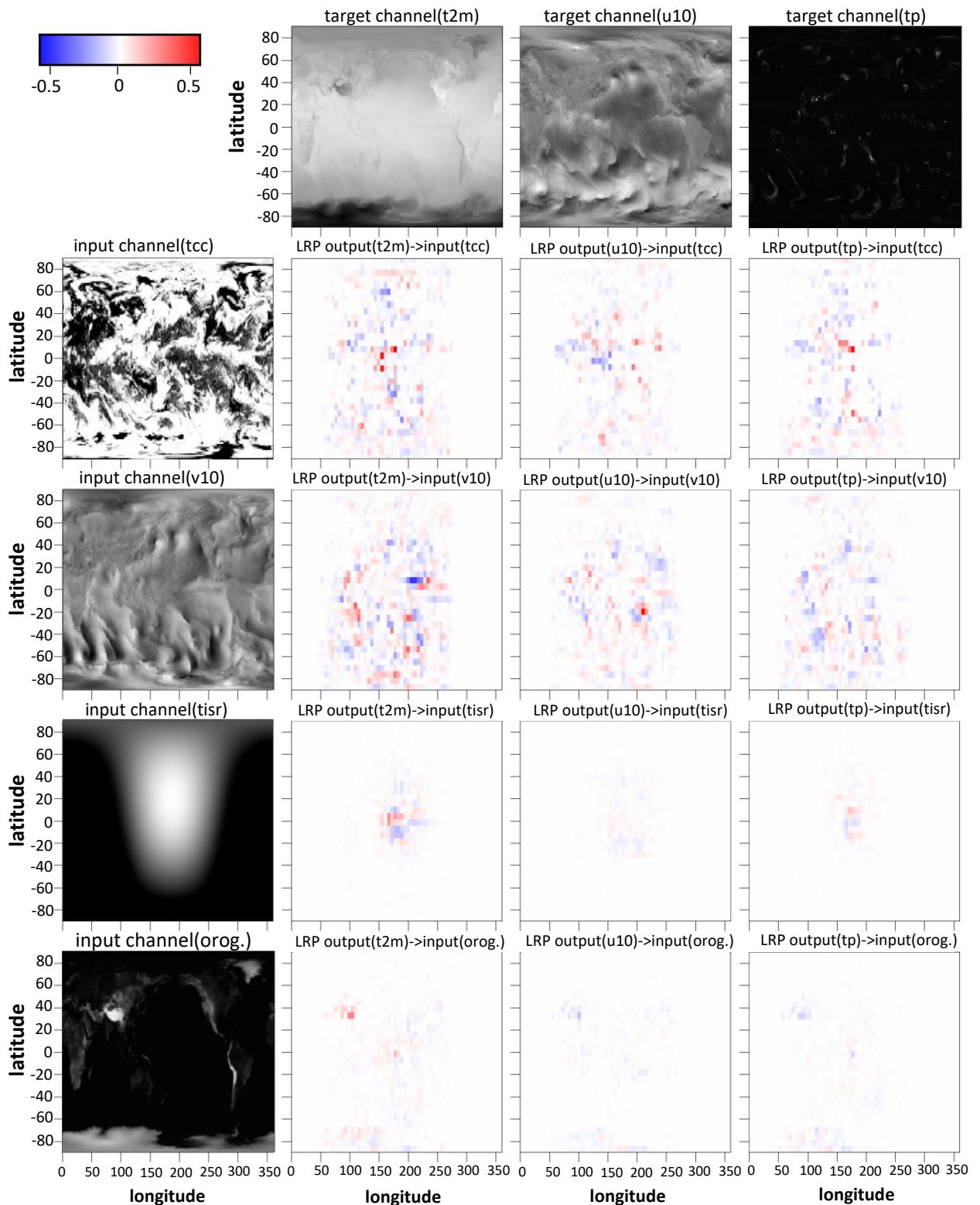
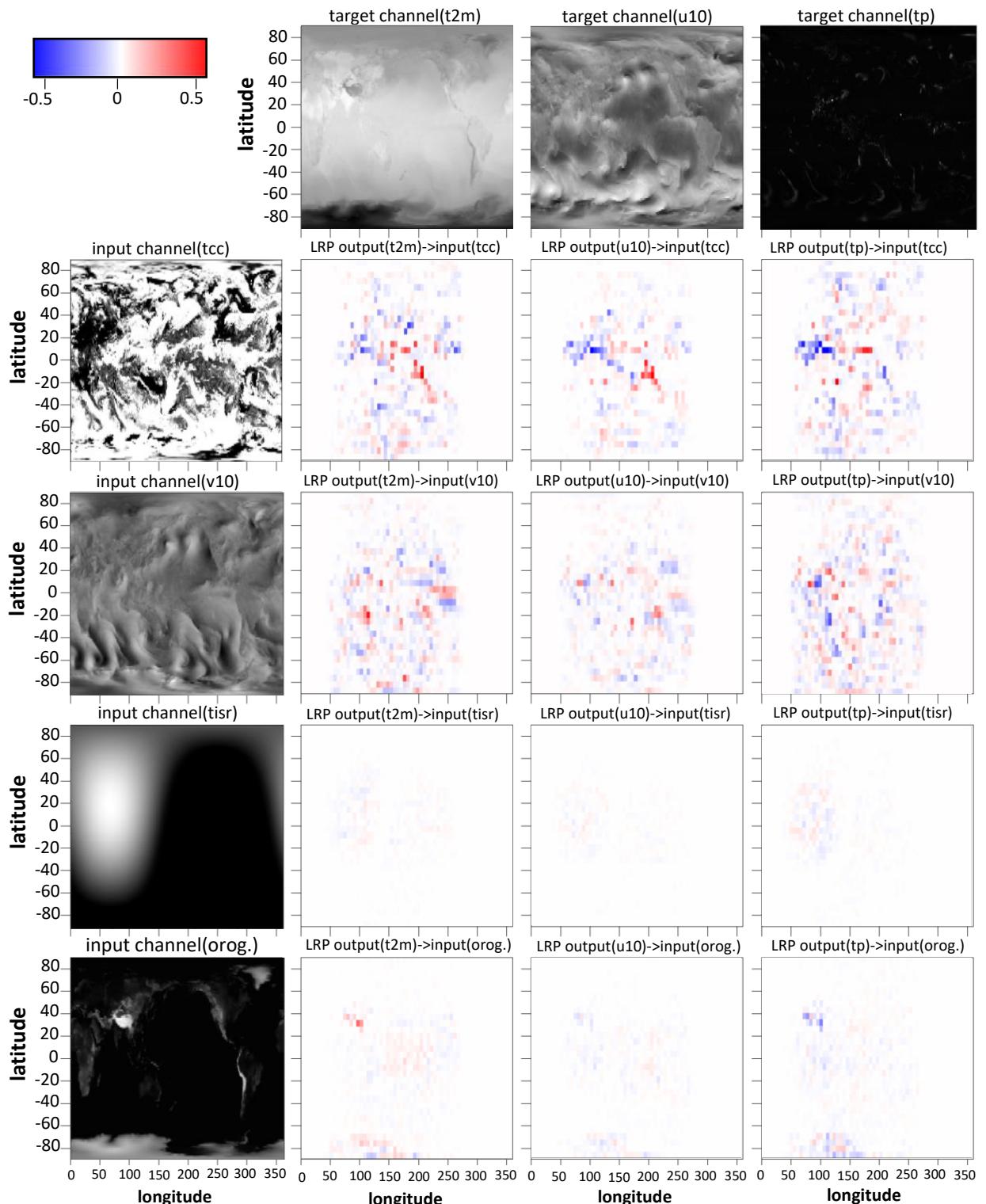


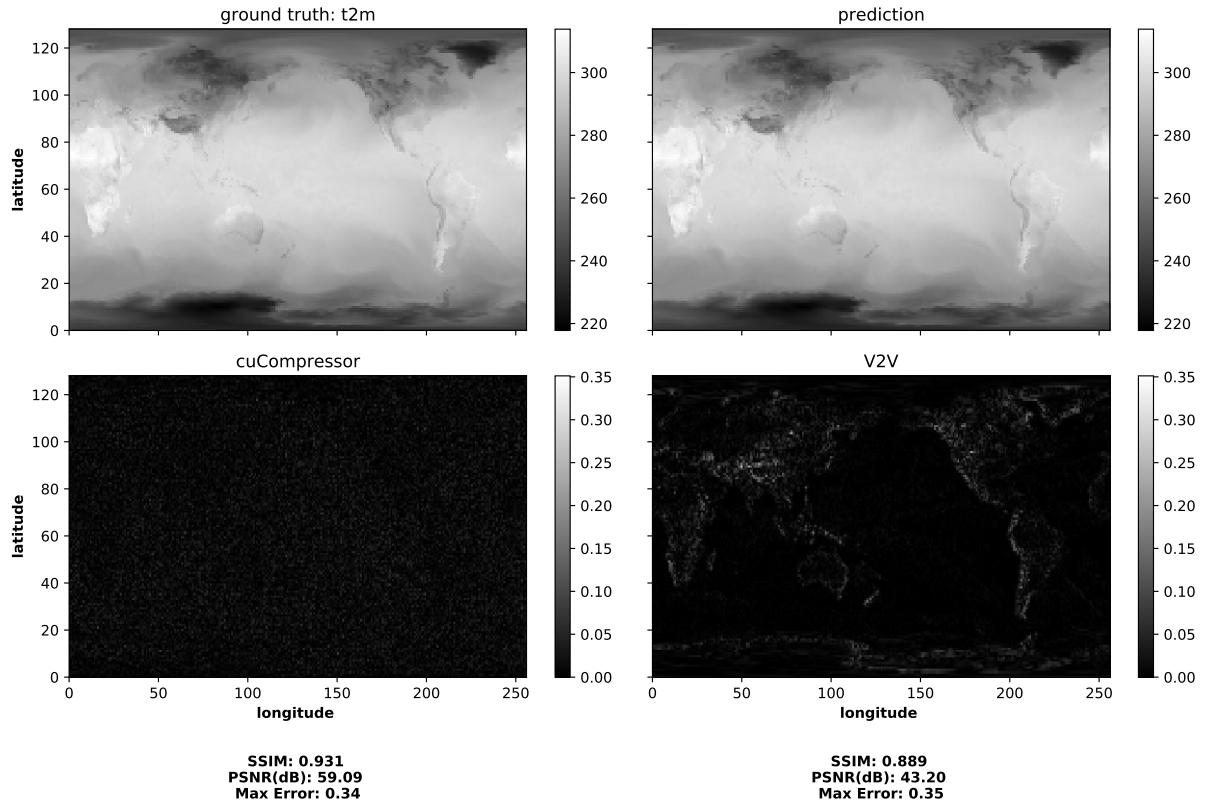
FIG. A2. Same as Fig. 7, but using UNet for training.



730 FIG. A3. LRP relevance maps with deviation-based selector function for the UNet model in the best input-output  
 731 configuration for WB data. Timestamp of data sample: May 15, 2004, 12h.



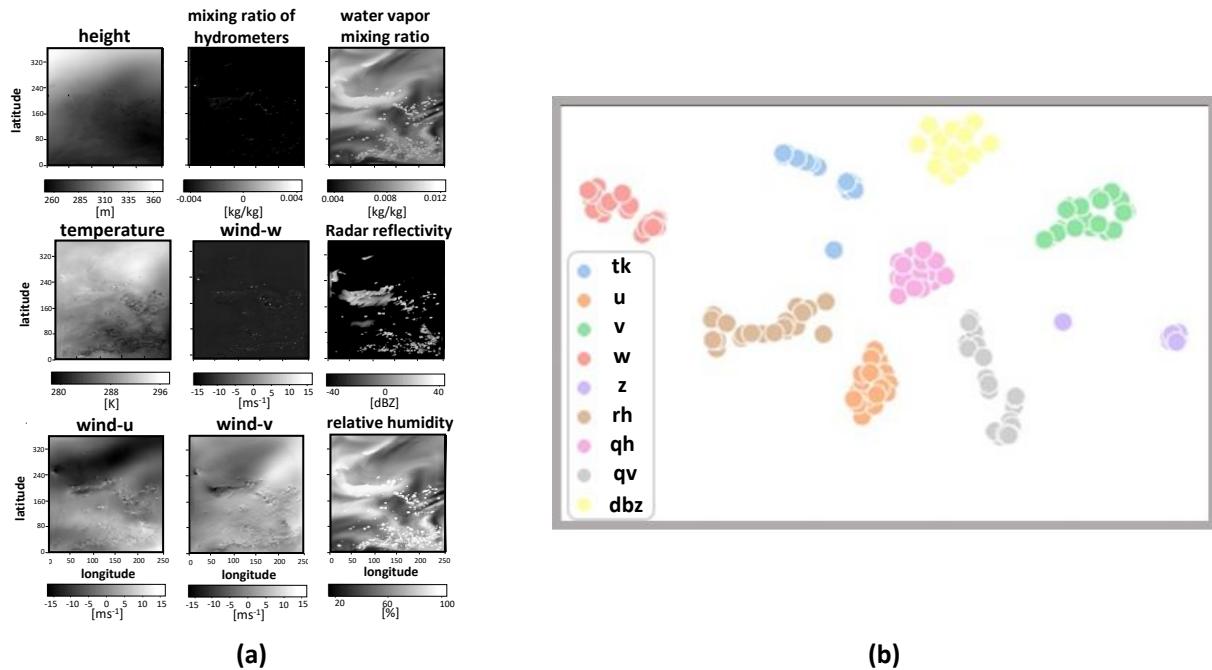
732 FIG. A4. LRP relevance maps with deviation-based selector function for the UNet model in the best input-output  
 733 configuration for WB data. Timestamp of data sample: May 15, 2004, 20h.



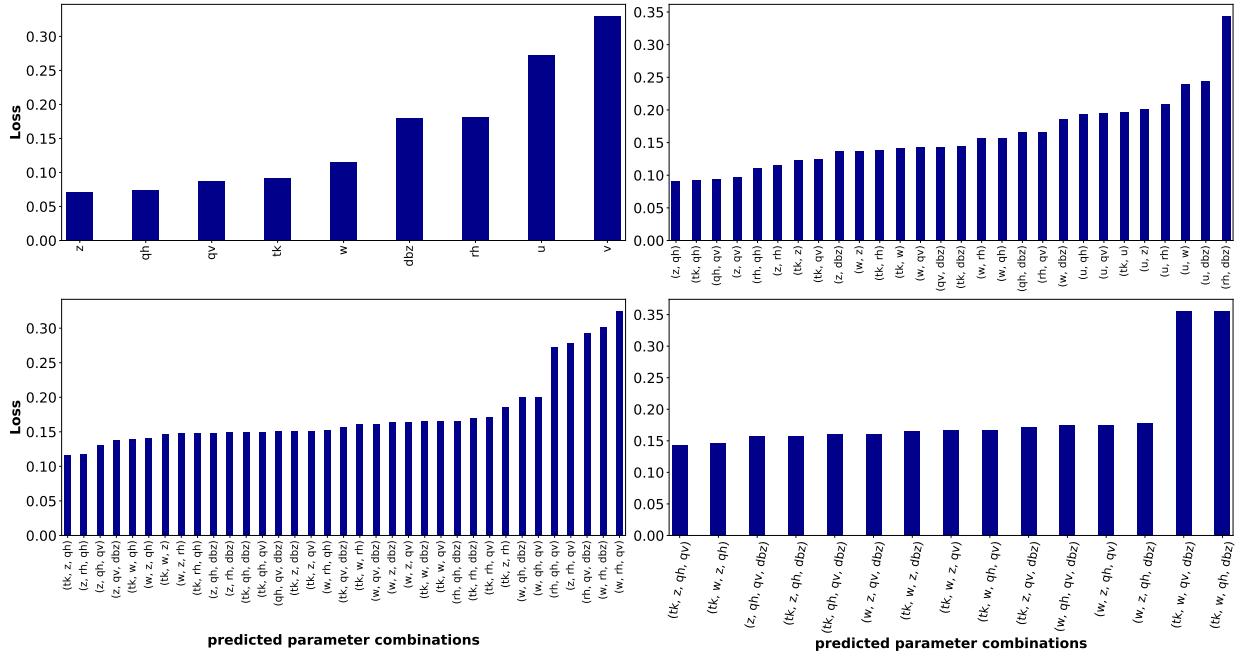
734 FIG. A5. Quality comparison of V2V against a dedicated compression algorithm for volumetric data. Parameter  
 735 field t2m compressed at a rate of 12:1 with the publicly available CUDA compression library by Treib et al.  
 736 (2012), which provides lossy compression using a combination of the discrete wavelet transform, coefficient  
 737 quantization, run-length encoding, and Huffman coding. Top left: Original field, top right: Parameter field  
 738 predicted using 3-to-3 V2V transfer. Bottom: Pixel-wise differences for reconstructed compressed field and  
 739 V2V reconstruction.

## Variable-to-Variable Transfer for the Convective-Scale Ensemble

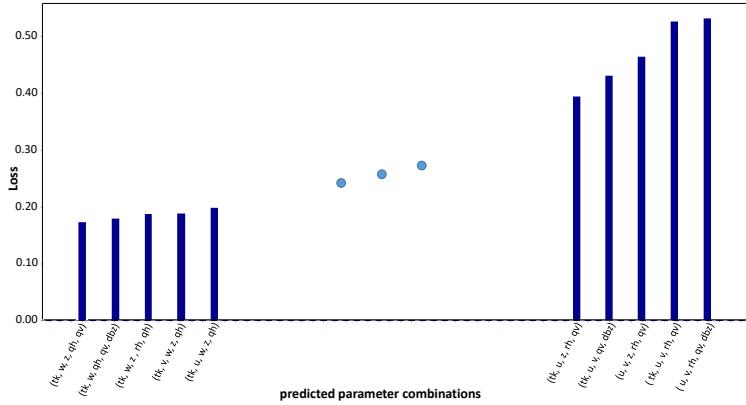
742 The appendix provides additional figures illustrating V2V transfer in the second dataset, the  
 743 convective-scale ensemble (CSEns) by Necker et al. (2020), which were excluded from the main  
 744 paper to improve readability.



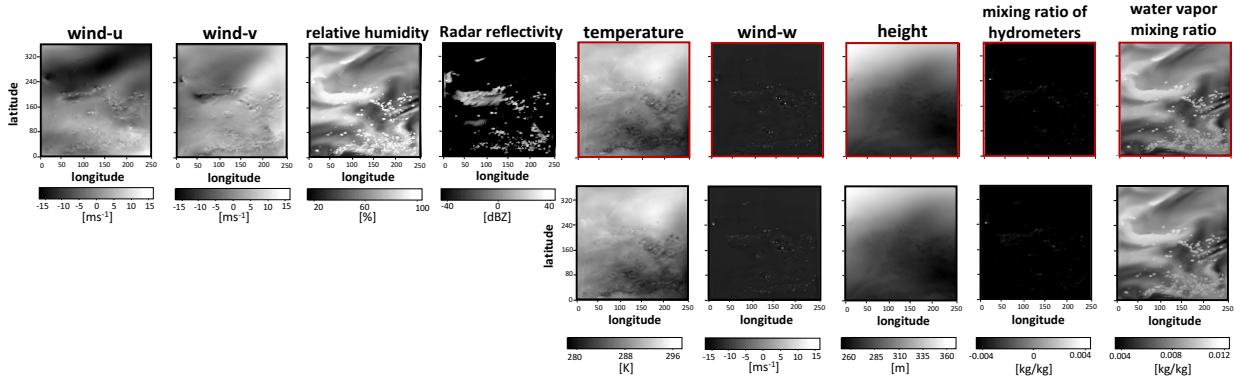
745 FIG. B1. Different parameter fields in the CSEns dataset. a) Gray-scale visualizations of the parameter fields  
 746 at a particular time. b) t-SNE projections of latent-space features of the parameter fields (different parameters  
 747 indicated by colors) at different times (note that projections for different initializations of t-SNE yield similar  
 748 groupings).



749 FIG. B2. Bar charts showing the losses of all networks trained for 4-to-5 parameter transfer with the CSEns  
 750 dataset using the proposed iterative loss-based approach. Top left: first iteration, top right: second iteration,  
 751 bottom left: third iteration, bottom right: fourth iteration. Bars represent losses after five epochs of training.



752 FIG. B3. Bar chart showing the losses of the best (left) and worst (right) possible networks for 4-to-5 parameter  
 753 transfer with the CSEns dataset. All  $\binom{9}{5}$  possible models have been trained for five epochs. Configurations with  
 754 intermediate losses have been omitted from the chart for clarity of the visualization.



755 FIG. B4. Reconstruction results for the CSEns datasets when the network is trained to predict five parameter  
 756 fields from four input fields. Top: The initial parameter fields. A red outline indicates those fields the network  
 757 has learned to predict from the others. Bottom: Predicted parameter fields.

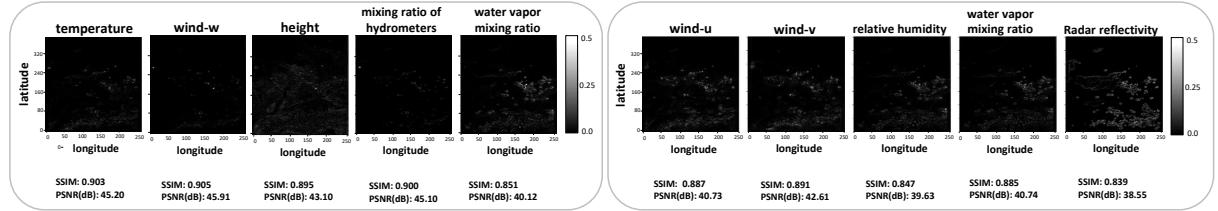


FIG. B5. Same as Fig. 8, but using the CSEns dataset.

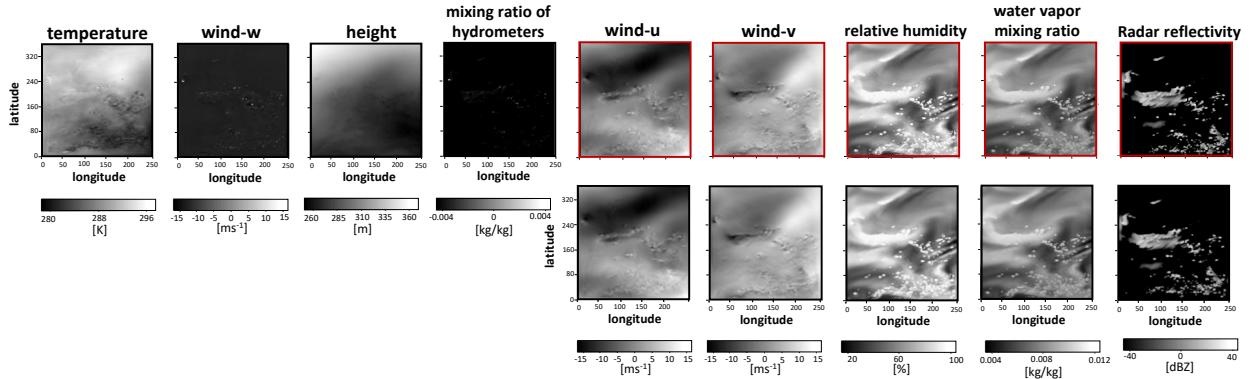


FIG. B6. Same as Fig. A1, but using the CSEns dataset.

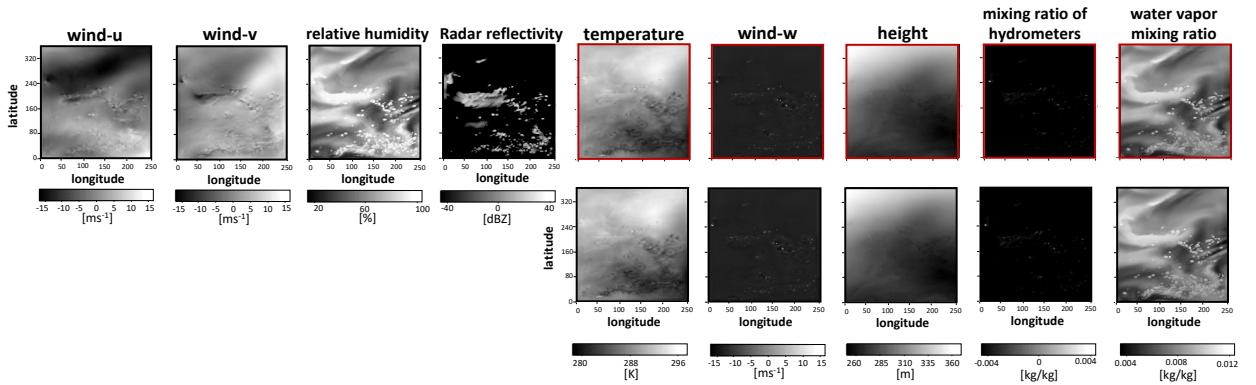
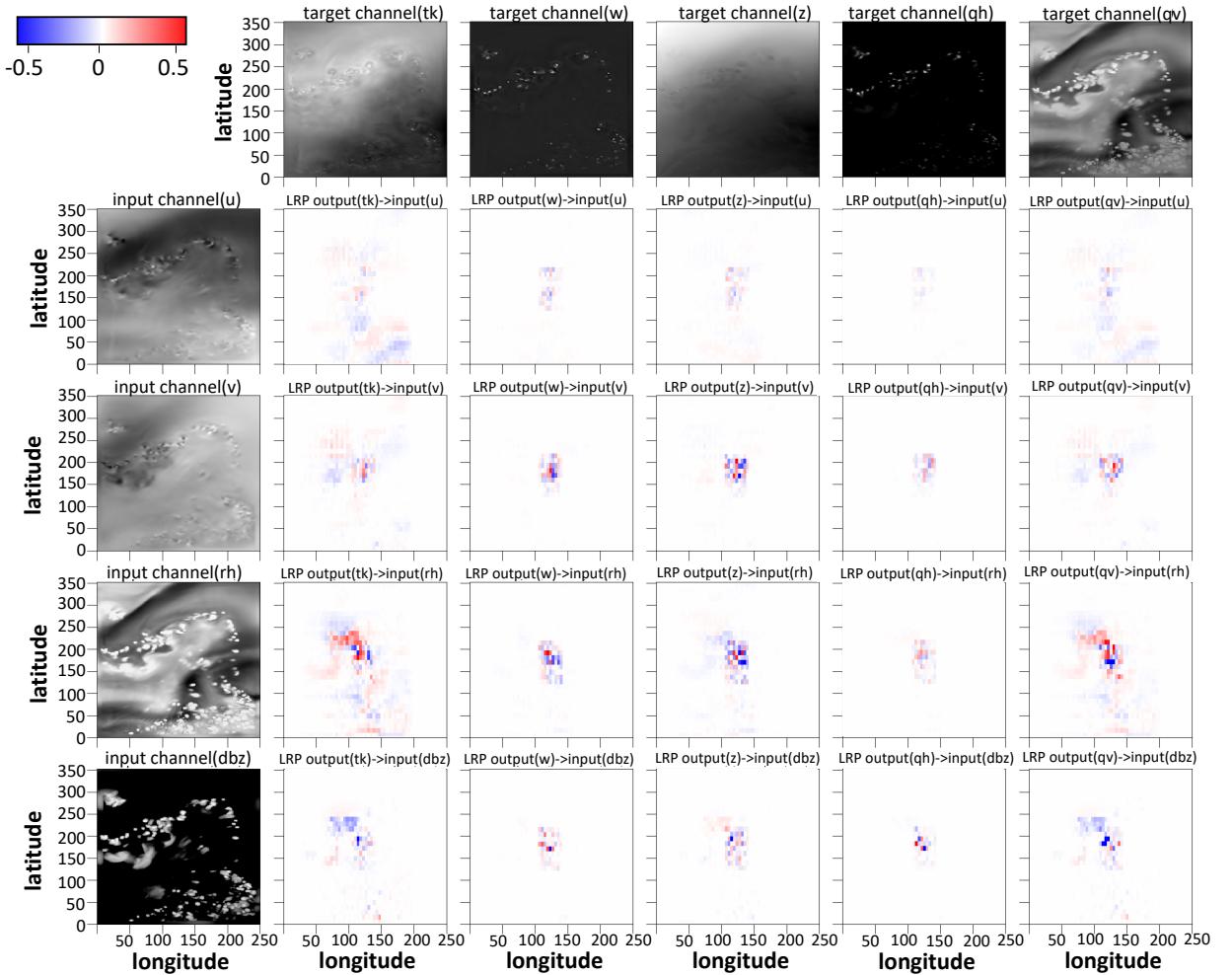
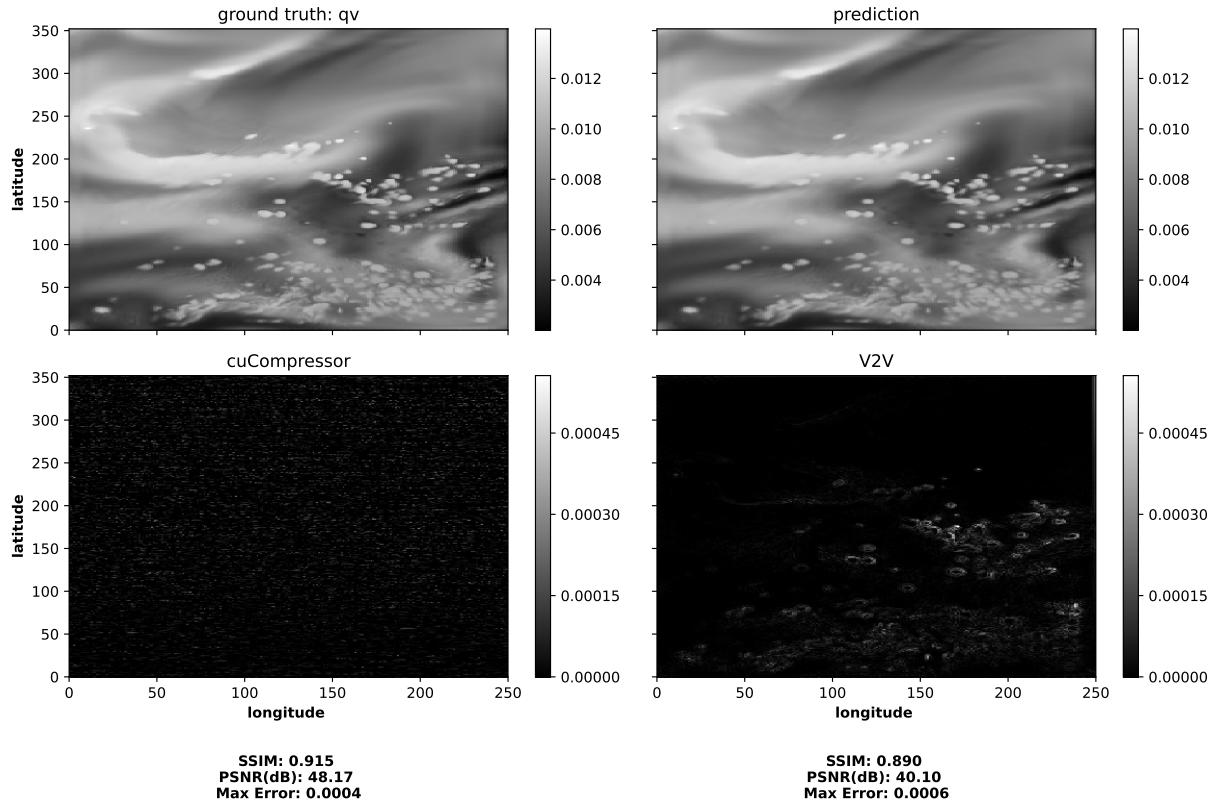


FIG. B7. Same as Fig. B4, but using the UNet architecture instead of the ResNet.



758 FIG. B8. LRP relevance maps with deviation-based selector function for the UNet model in the best input-output  
759 configuration for CSEns data. Timestamp of data sample: June 1, 2016, 17h.



760 FIG. B9. Quality comparison of V2V against a dedicated compression algorithm for volumetric data. Parameter  
 761 field qv compressed at a rate of 12:1 with the publicly available CUDA compression library by Treib et al. (2012),  
 762 which provides lossy compression using a combination of the discrete wavelet transform, coefficient quantization,  
 763 run-length encoding, and Huffman coding. Top left: Original field, top right: Parameter field predicted using  
 764 4-to-5 V2V transfer. Bottom: Pixel-wise differences for reconstructed compressed field and V2V reconstruction.