

# AI for High-Resolution Regional Climate Modeling

Advanced School on  
High-Performance Computing and  
Applied AI for High-Resolution Regional  
Climate Modeling



The Abdus Salam  
**International Centre  
for Theoretical Physics**



Jose González-Abad  
Postdoctoral Researcher  
[gonzabad@ifca.unican.es](mailto:gonzabad@ifca.unican.es)



**i F C A**  
Instituto de Física de Cantabria



**CSIC**  
CONSEJO SUPERIOR DE INVESTIGACIONES CIENTÍFICAS



**UC** | Universidad  
de Cantabria

# Table of Contents

## **Part I: Review**

- DL for Weather/Climate
- RCM Emulation
- Conclusions

## **Part II: Training a RCM emulator**

- Problem Statement
- Training Frameworks
- Training and Inference

## **Part III: Evaluating a RCM emulator**

- Importance of Proper Evaluation
- Soft-Transferability
- Hard-Transferability
- Metrics

## **Part IV: The Importance of Benchmarks**

- What is a benchmark?
- CORDEX-ML-Bench

# Table of Contents

## **Part I: Review**

- DL for Weather/Climate
- RCM Emulation
- Conclusions

## **Part II: Training a RCM emulator**

- Problem Statement
- Training Frameworks
- Training and Inference

## **Part III: Evaluating a RCM emulator**

- Importance of Proper Evaluation
- Soft-Transferability
- Hard-Transferability
- Metrics

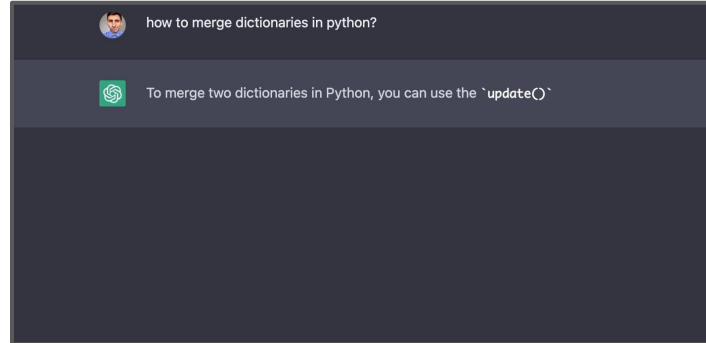
## **Part IV: The Importance of Benchmarks**

- What is a benchmark?
- CORDEX-ML-Bench

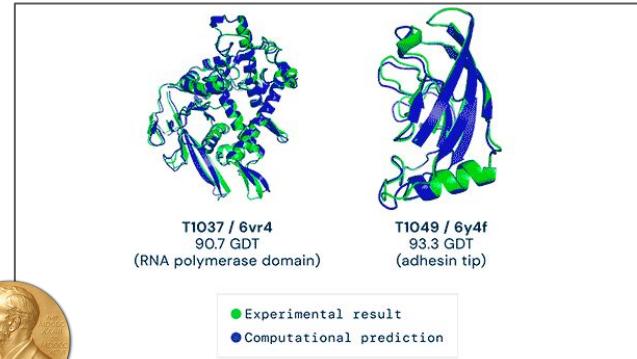
# Deep Learning Revolution

Deep learning models have **revolutionized many fields**, achieving results that were never imagined a few years ago.

LLMs



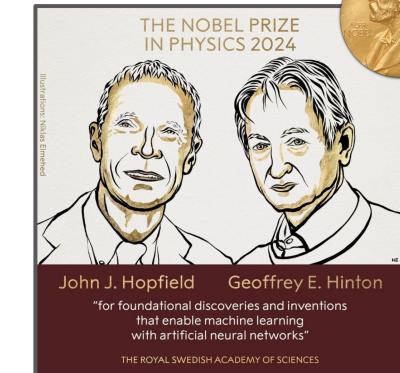
Protein Folding



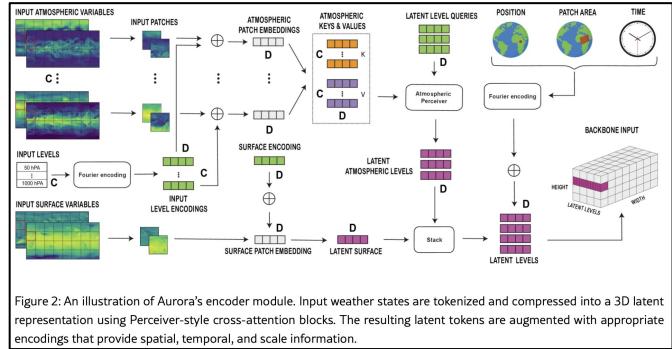
Video Generation



Mastering of Go



# Deep Learning Revolution

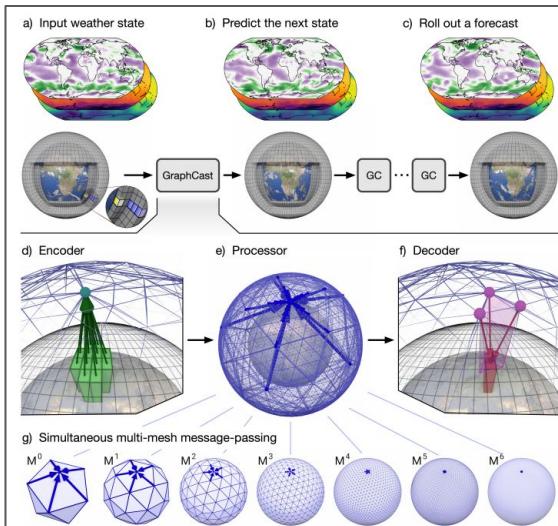


## A Foundation Model for the Earth System

Cristian Bodnar<sup>1,2†</sup>, Wessel P. Bruinsma<sup>1†</sup>, Ana Lucic<sup>1,3†</sup>, Megan Stanley<sup>1†</sup>, Anna Vaughan<sup>1†</sup>, Johannes Brandstetter<sup>1,5</sup>, Patrick Garvan<sup>1</sup>, Maik Riechert<sup>1</sup>, Jonathan A. Weyn<sup>6</sup>, Haiyu Dong<sup>6</sup>, Jayesh K. Gupta<sup>2,7</sup>, Kit Thambiratnam<sup>6</sup>, Alexander T. Archibald<sup>4</sup>, Chun-Chieh Wu<sup>8</sup>, Elizabeth Heider<sup>1</sup>, Max Welling<sup>1,3</sup>, Richard E. Turner<sup>1,4,9</sup>, Paris Perdikaris<sup>1,10\*</sup>

## Spherical Fourier Neural Operators: Learning Stable Dynamics on the Sphere

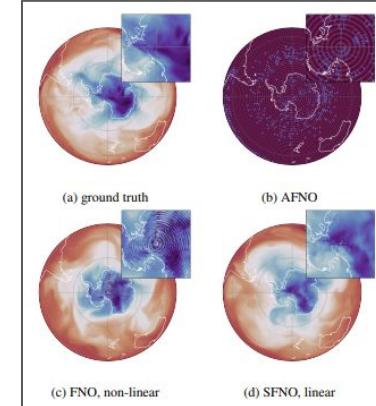
Boris Bonev<sup>1</sup> Thorsten Kurth<sup>1</sup> Christian Hundt<sup>1</sup> Jaideep Pathak<sup>1</sup> Maximilian Baust<sup>1</sup> Karthik Kashinath<sup>1</sup>  
Anima Anandkumar<sup>1,2</sup>



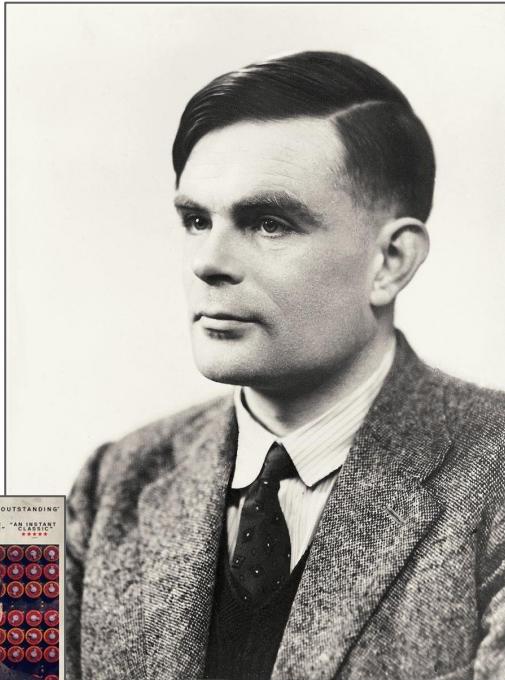
## GraphCast: Learning skillful medium-range global weather forecasting

Remi Lam<sup>\*1</sup>, Alvaro Sanchez-Gonzalez<sup>\*1</sup>, Matthew Willison<sup>\*1</sup>, Peter Wirsberger<sup>\*1</sup>, Meire Fortunato<sup>\*1</sup>, Ferran Alet<sup>\*1</sup>, Suman Ravuri<sup>\*1</sup>, Timo Ewalds<sup>1</sup>, Zach Eaton-Rosen<sup>1</sup>, Weihua Hu<sup>1</sup>, Alexander Merose<sup>2</sup>, Stephan Hoyer<sup>2</sup>, George Holland<sup>1</sup>, Oriol Vinyals<sup>1</sup>, Jacklynn Stott<sup>1</sup>, Alexander Pritzel<sup>1</sup>, Shakir Mohamed<sup>1</sup> and Peter Battaglia<sup>1</sup>

\*equal contribution, <sup>1</sup>Google DeepMind, <sup>2</sup>Google Research



## Weather Turing Test



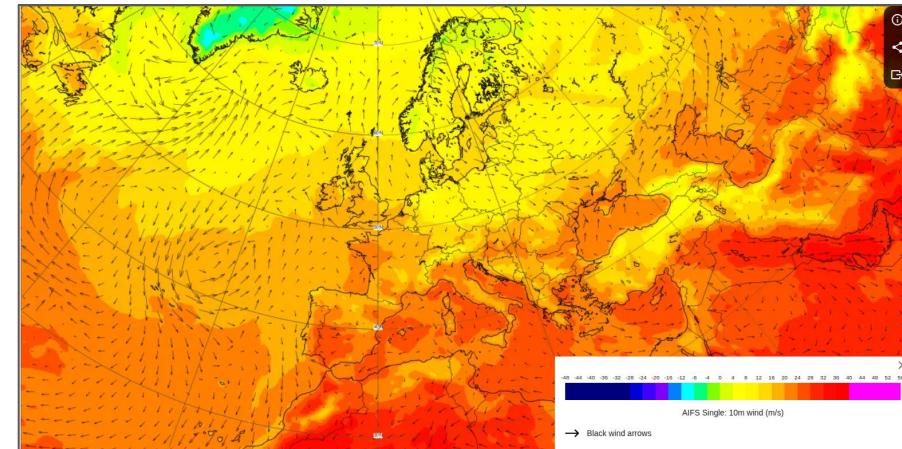
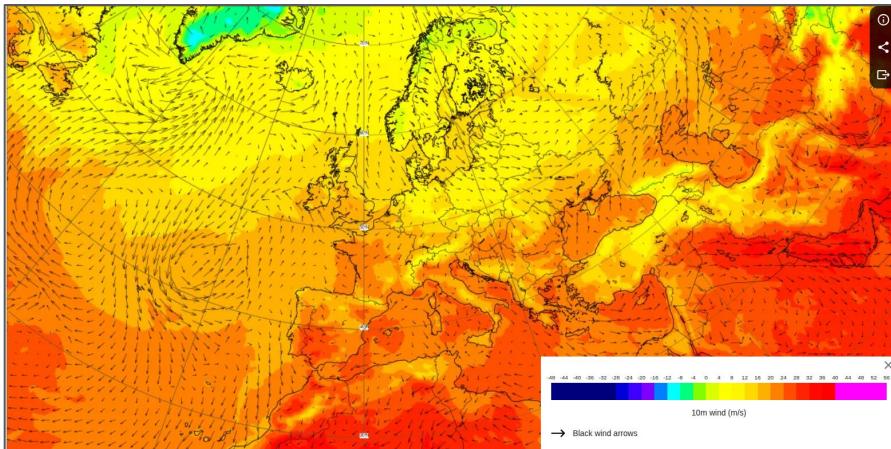
Alan Turing

The **Turing Test**, proposed by Alan Turing in 1950, is a benchmark for **determining whether a machine can exhibit intelligent behavior indistinguishable from that of a human**. In the test, a human evaluator interacts with both a machine and a human without knowing which is which. **If the evaluator cannot reliably tell them apart, the machine is said to have passed the test.**

# Weather Turing Test

Which one comes from a **physics-based** model and which from a **DL-based** one?

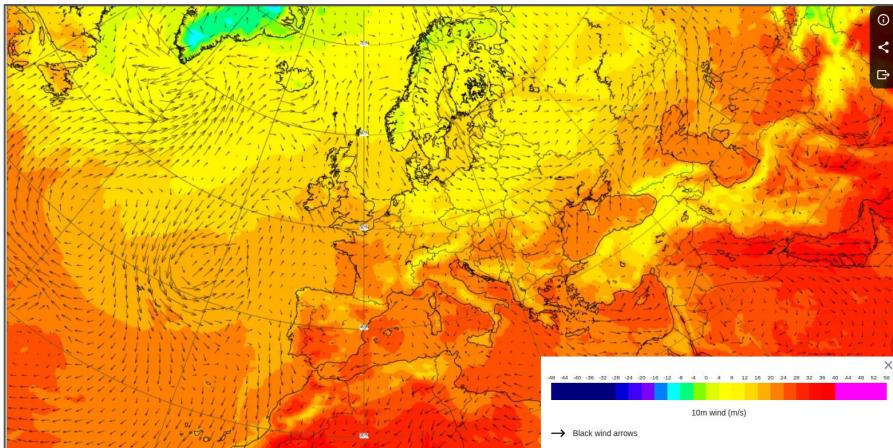
**2m temperature and 10m wind**



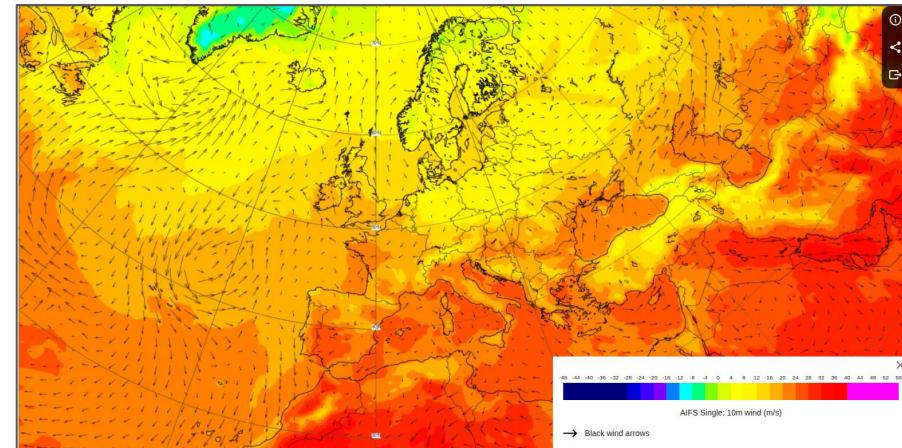
# Weather Turing Test

Which one comes from a **physics-based** model and which from a **DL-based** one?

**2m temperature and 10m wind**



**Physics-based model**

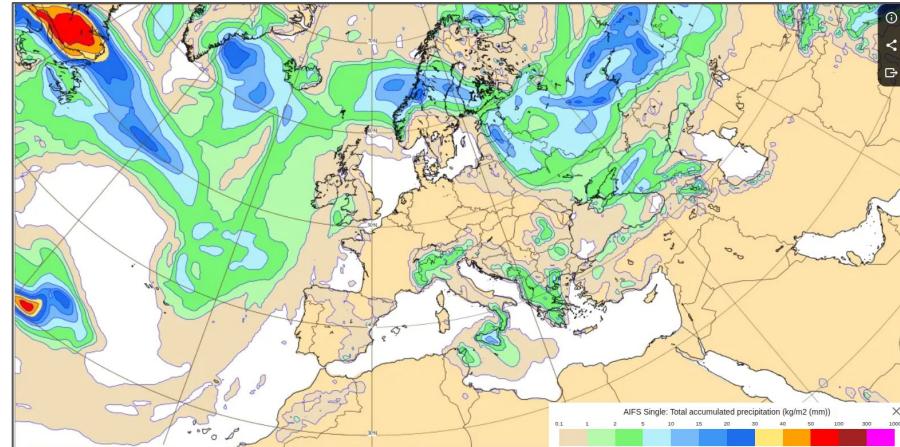
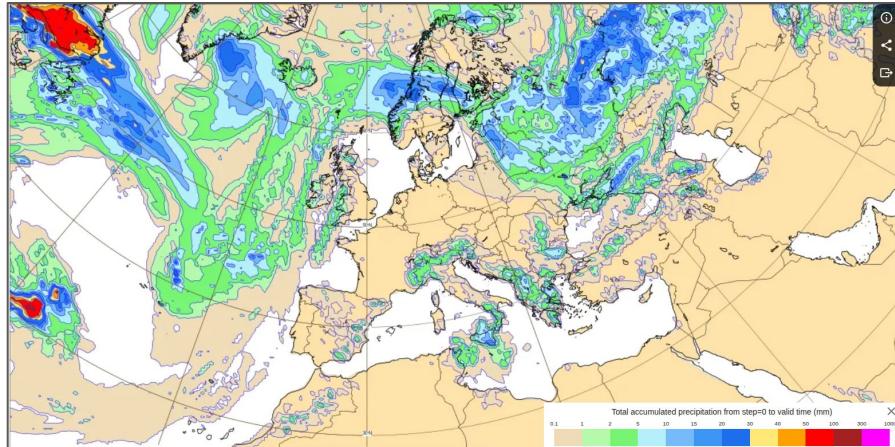


**AIFS (DL-based model)**

# Weather Turing Test

Which one comes from a **physics-based** model and which from a **DL-based** one?

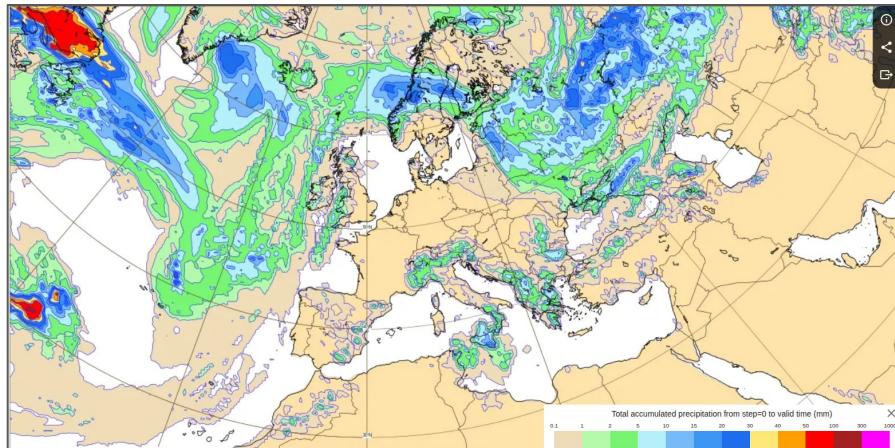
**Accumulated Precipitation**



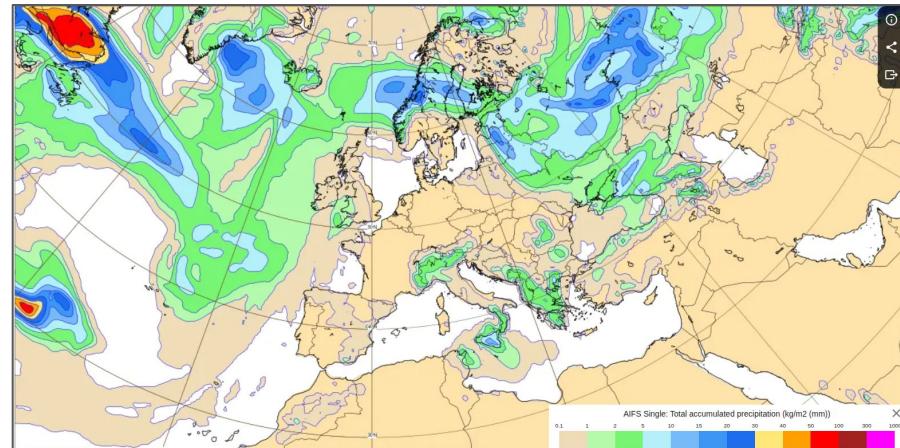
# Weather Turing Test

Which one comes from a **physics-based** model and which from a **DL-based** one?

## Accumulated Precipitation



Physics-based model

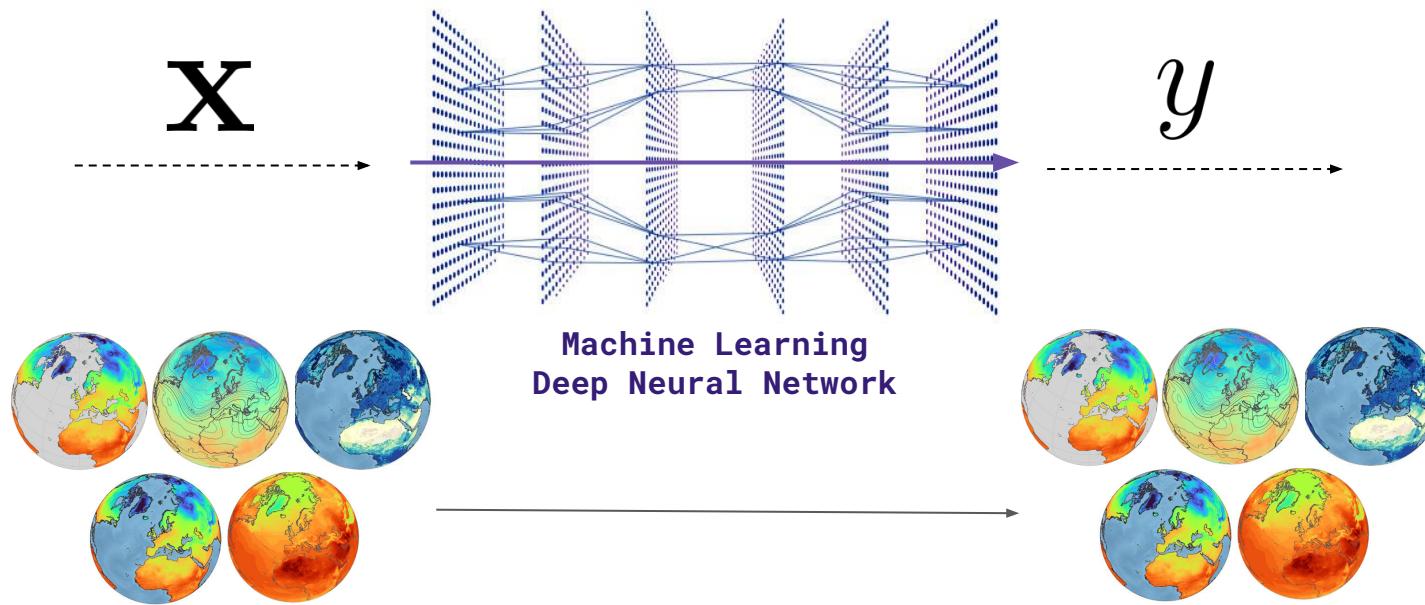


AIFS (DL-based model)

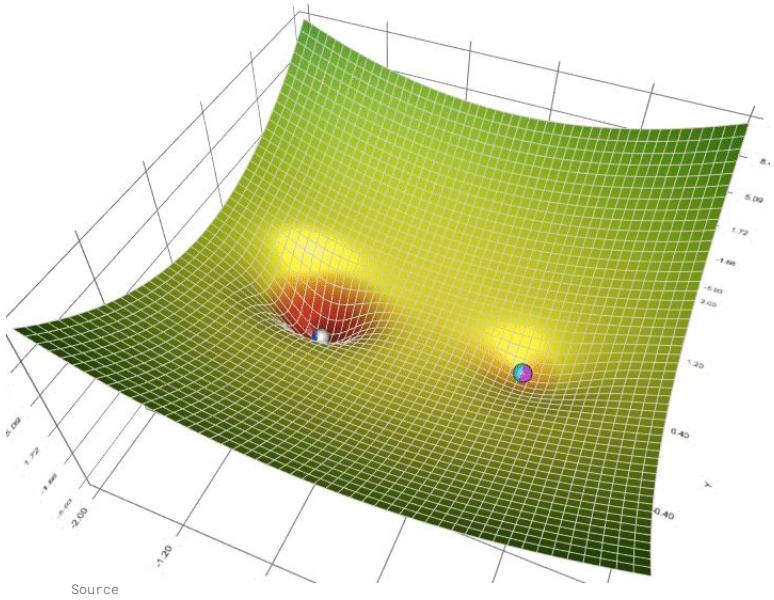
## Weather Turing Test

**DL-based** models tend to generate **smoothed/blurry predictions**

# Deep Learning



# Deep Learning



We adjust the parameters of the neural network to minimize the value of a loss function associated with the task we want to accomplish

The most typical is the **Mean Squared Error** (MSE)

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

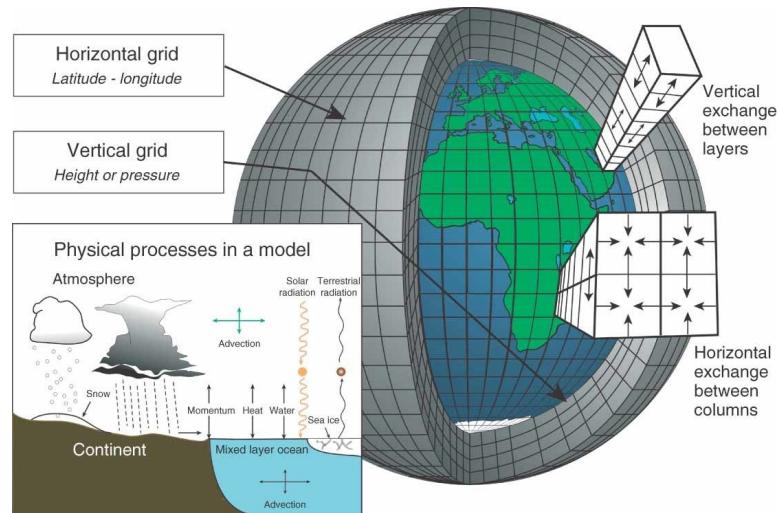
# Global Climate Models

**Global Climate Models** (GCM) are the main tools used to **simulate the evolution of climate** by numerically solving the physical equations characterizing climate dynamics:

## Climate equations

$$\begin{aligned}\frac{d\mathbf{v}}{dt} &= -\alpha \nabla p - \nabla \phi + \mathbf{F} - 2\Omega \times \mathbf{v} \\ \frac{\partial \rho}{\partial t} &= -\nabla \cdot (\rho \mathbf{v}) \\ p\alpha &= RT \\ Q &= C_p \frac{dT}{dt} - \alpha \frac{dp}{dt} \\ \frac{\partial \rho q}{\partial t} &= -\nabla \cdot (\rho \mathbf{v} q) + \rho(E - C)\end{aligned}$$

discretize over space and time



Source: Edwards, P. N. (2011). History of climate modeling. Wiley Interdisciplinary Reviews: Climate Change, 2(1), 128-139.

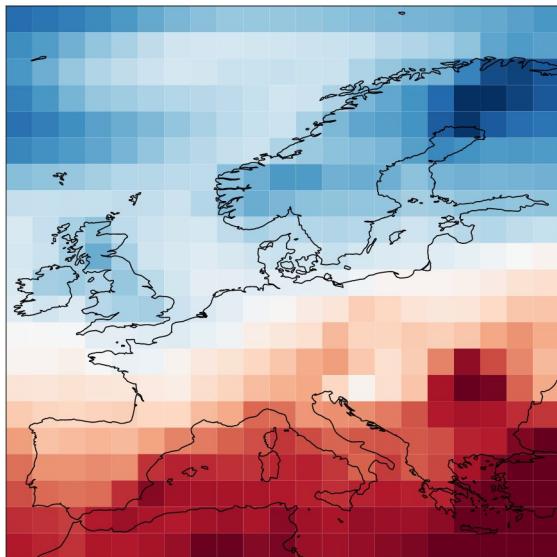
With **typical resolutions** of hundreds of km (e.g. **~100 km** in CMIP6)

## Coarse Resolution of GCMs

Due to deficiencies in the climate modeling process, **GCMs can not accurately reproduce the local scale**

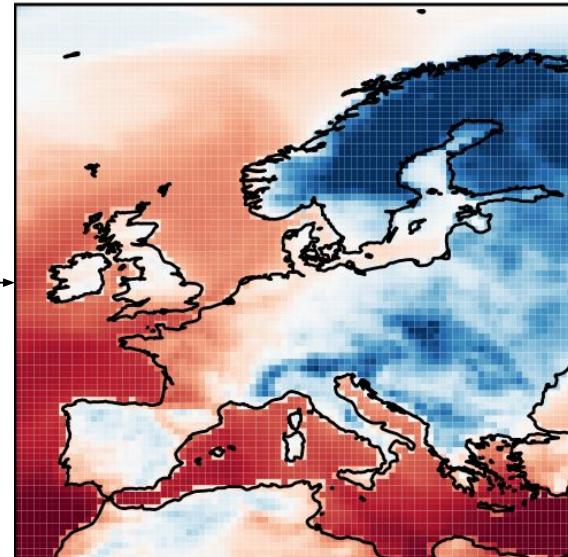
This makes it difficult to use GCMs in different socio-economical activities where **regional information** is required

GCM output



How can we avoid this difference in resolution?

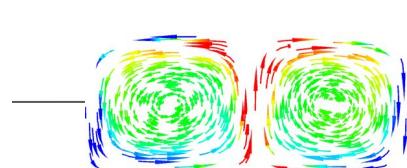
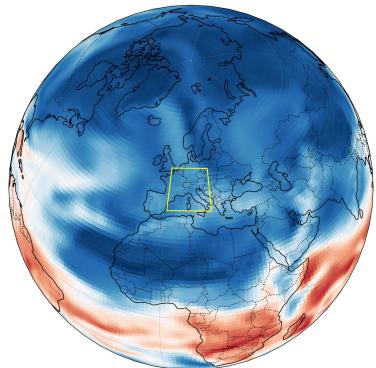
Observations



# Regional Climate Models

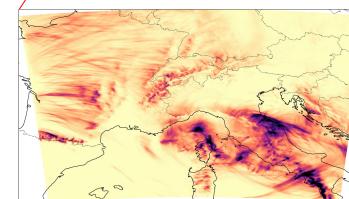
**Global Climate Model (GCM)**

CNRM-CM5 (~100 km)



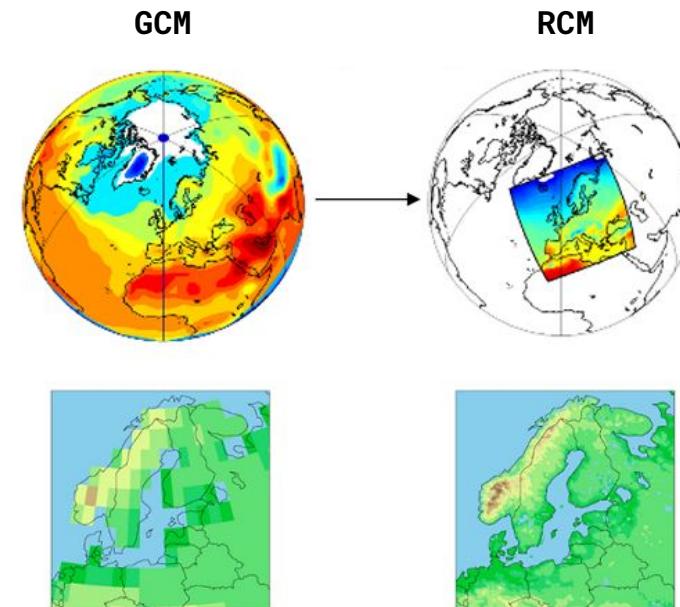
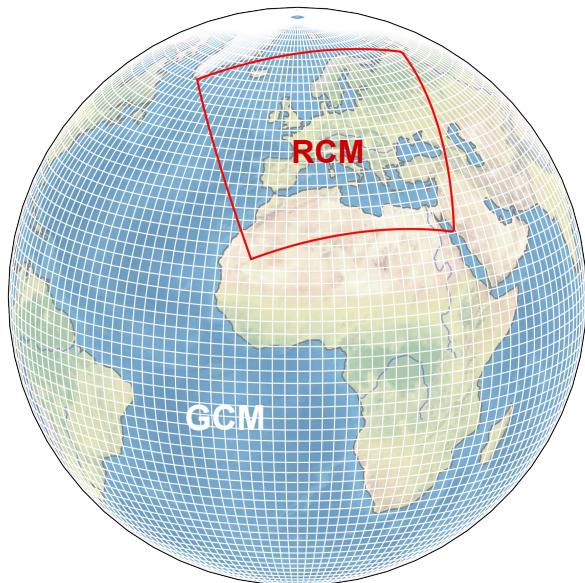
**Regional Climate Model (RCM)**

AROME (~3 km)



A RCM is a **high-resolution physics-based model focused on a specific region**, using boundary conditions from a GCM.

## Regional Climate Models

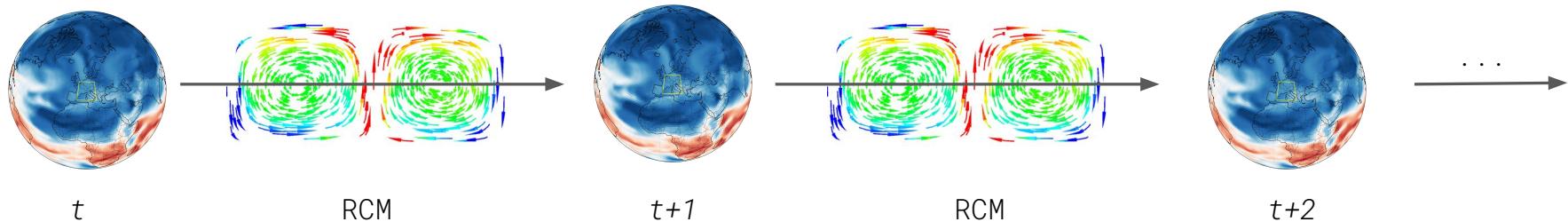


A RCM is a **high-resolution physics-based model focused on a specific region**, using boundary conditions from a GCM.

Source

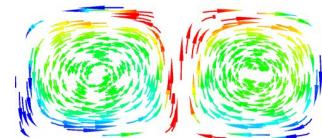
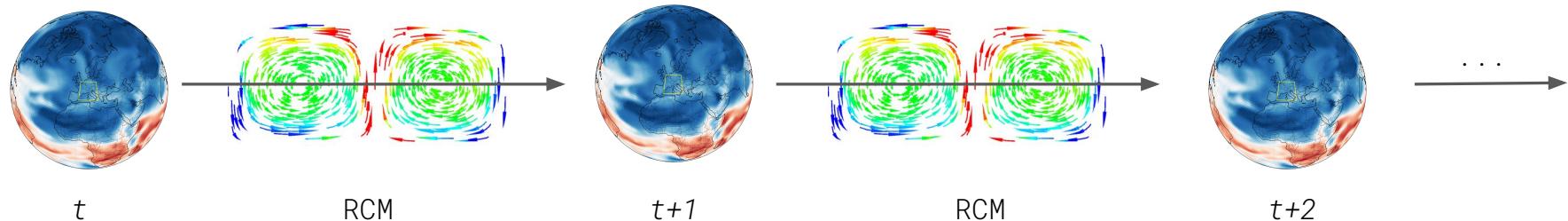
# Regional Climate Models

Running a RCM is **computationally expensive** and **challenging**

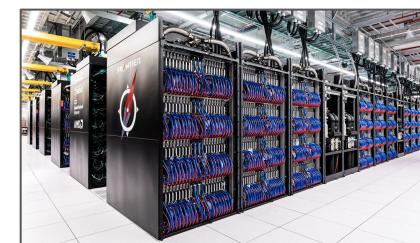


# Regional Climate Models

Running a RCM is **computationally expensive** and **challenging**

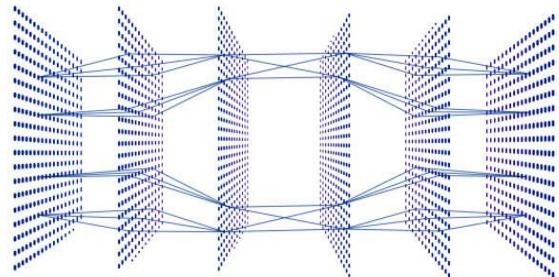


requires

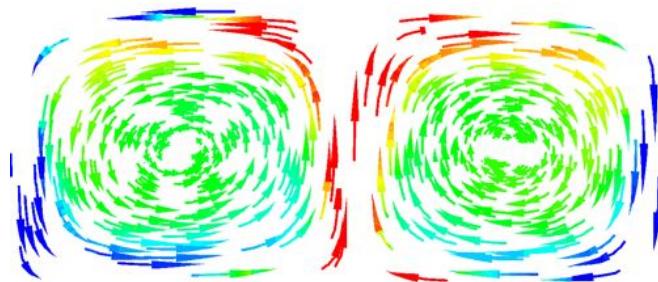


## RCM Emulation

What if we **emulate** the RCM with a Deep Learning model?



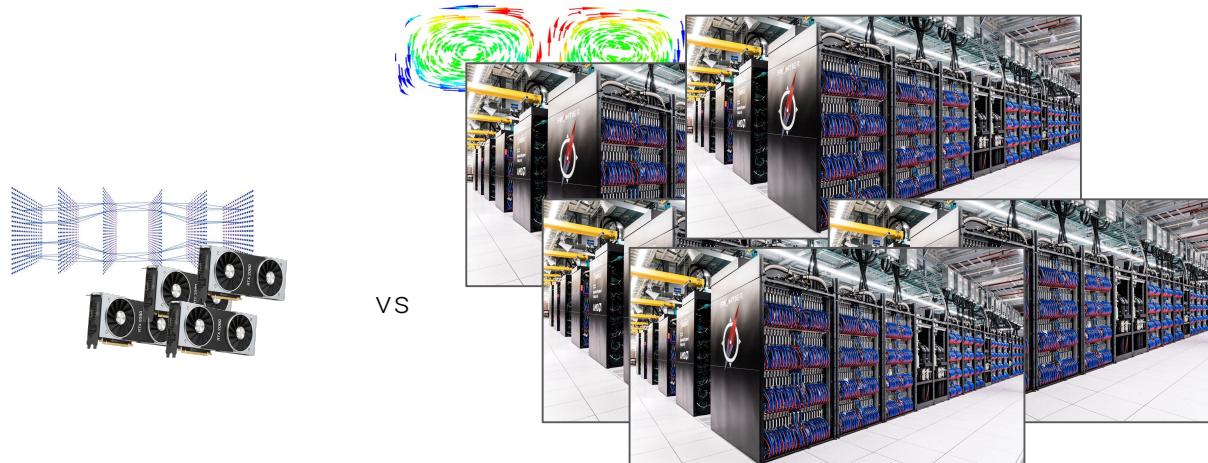
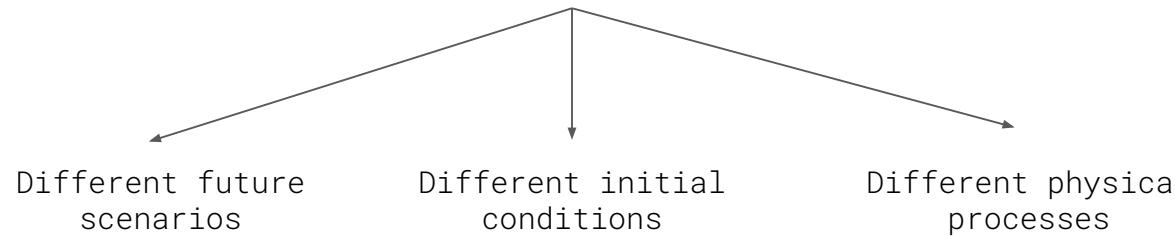
Deep Learning



RCM

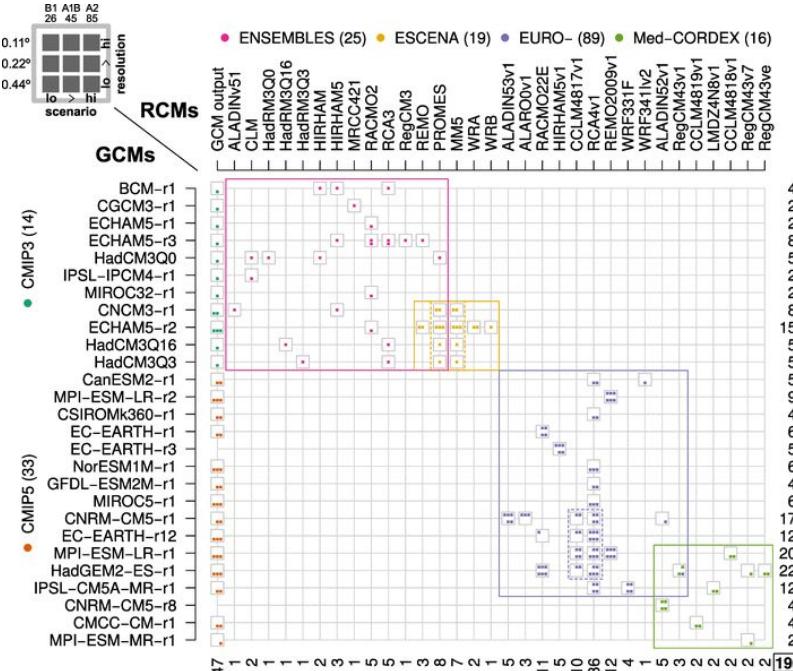
## Benefits of RCM Emulation

The main point of the emulation of RCMs is to better **assess the uncertainty linked to future projections**



# Benefits of RCM Emulation

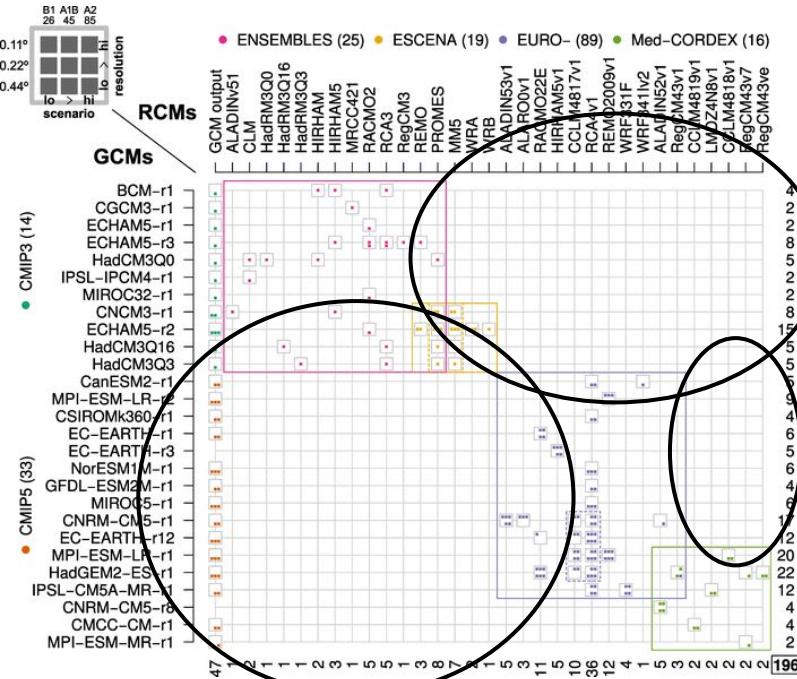
These emulators can help us **fill the GCM/RCM matrix**



Fuente: Fernández, J., Frias, M. D., Cabos, W. D., Cofiño, A. S., Dominguez, M., Fita, L., ... & Sánchez, E. (2019). Consistency of climate change projections from multiple global and regional model intercomparison projects. *Climate dynamics*, 52, 1139-1156.

# Benefits of RCM Emulation

These emulators can help us **fill the GCM/RCM matrix**



Fuente: Fernández, J., Frias, M. D., Cabos, W. D., Cofiño, A. S., Dominguez, M., Fita, L., ... & Sánchez, E. (2019). Consistency of climate change projections from multiple global and regional model intercomparison projects. Climate dynamics, 52, 1139-1156.

## How are RCM emulators trained?

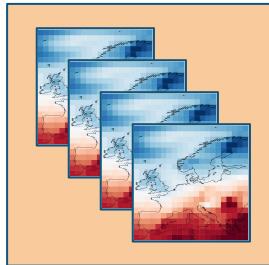
How can we emulate a RCM with DL? \*



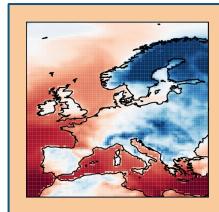
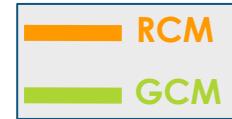
## How are RCM emulators trained?

How can we emulate a RCM with DL? \*

Large-scale variables  
(Low-resolution)



*Air-temperature  
Specific Humidity  
Wind Components  
Geopotential*



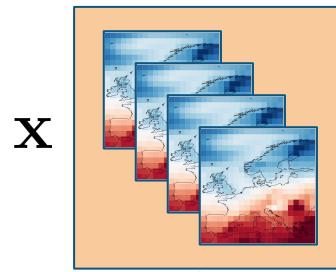
*Temperature  
Precipitation*

Surface variable  
(High-resolution)

## How are RCM emulators trained?

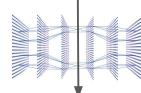
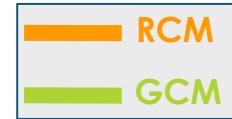
How can we emulate a RCM with DL? \*

Large-scale variables  
(Low-resolution)

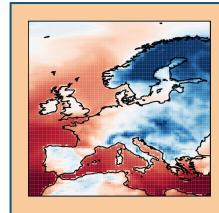


$X$

Air-temperature  
Specific Humidity  
Wind Components  
Geopotential



$y$



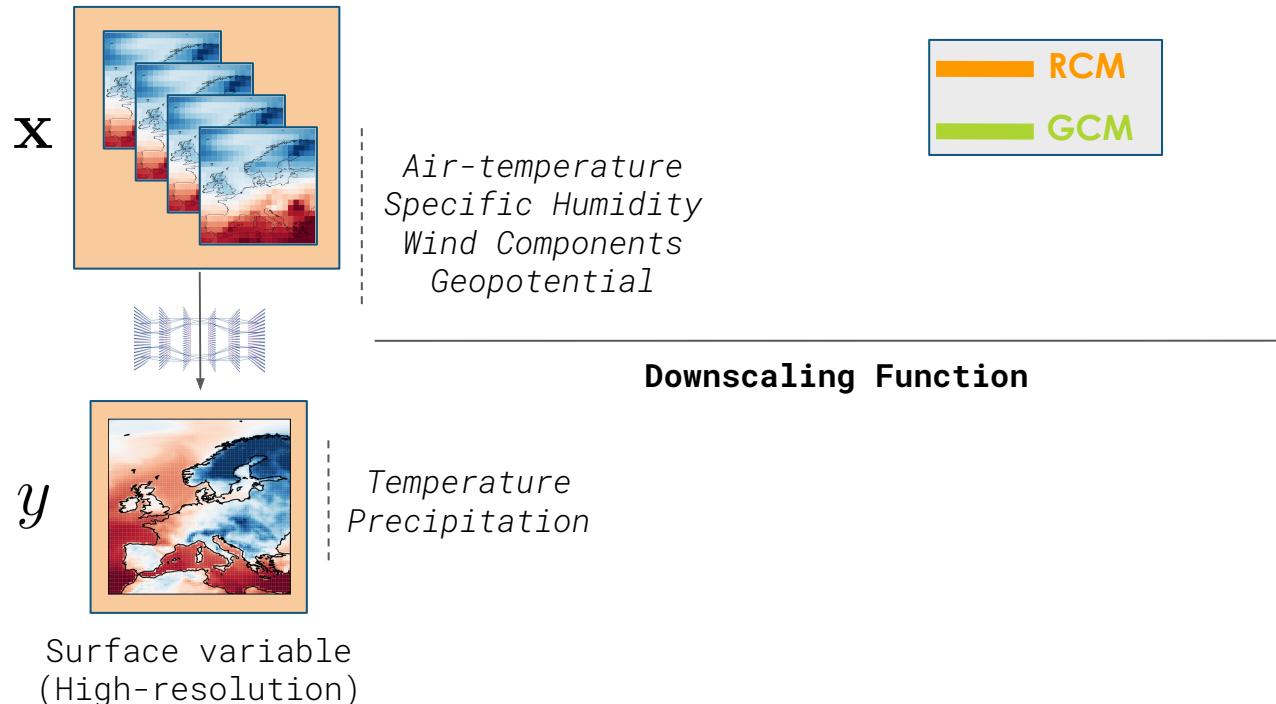
Temperature  
Precipitation

Surface variable  
(High-resolution)

## How are RCM emulators trained?

How can we emulate a RCM with DL? \*

Large-scale variables  
(Low-resolution)

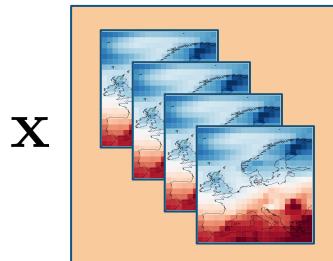


\* RCM emulation generally focuses on local-scale variables (e.g., temperature and precipitation)

## How are RCM emulators trained?

How can we emulate a RCM with DL? \*

Large-scale variables  
(Low-resolution)

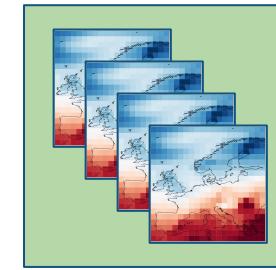


$X$

Air-temperature  
Specific Humidity  
Wind Components  
Geopotential



Large-scale variables  
(Low-resolution)



$y$

Temperature  
Precipitation

Downscaling Function

Surface variable  
(High-resolution)

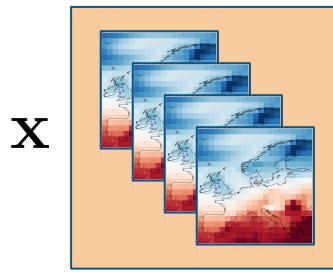
\*

RCM emulation generally focuses on local-scale variables (e.g., temperature and precipitation)

## How are RCM emulators trained?

How can we emulate a RCM with DL? \*

Large-scale variables  
(Low-resolution)

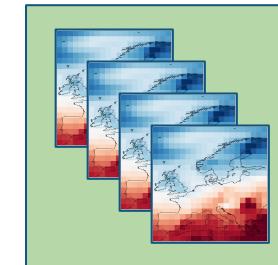


$X$

Air-temperature  
Specific Humidity  
Wind Components  
Geopotential



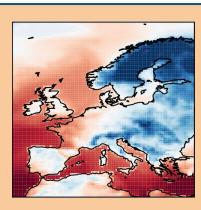
Large-scale variables  
(Low-resolution)



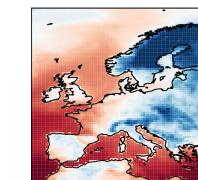
$y$

Temperature  
Precipitation

Surface variable  
(High-resolution)



Downscaling Function



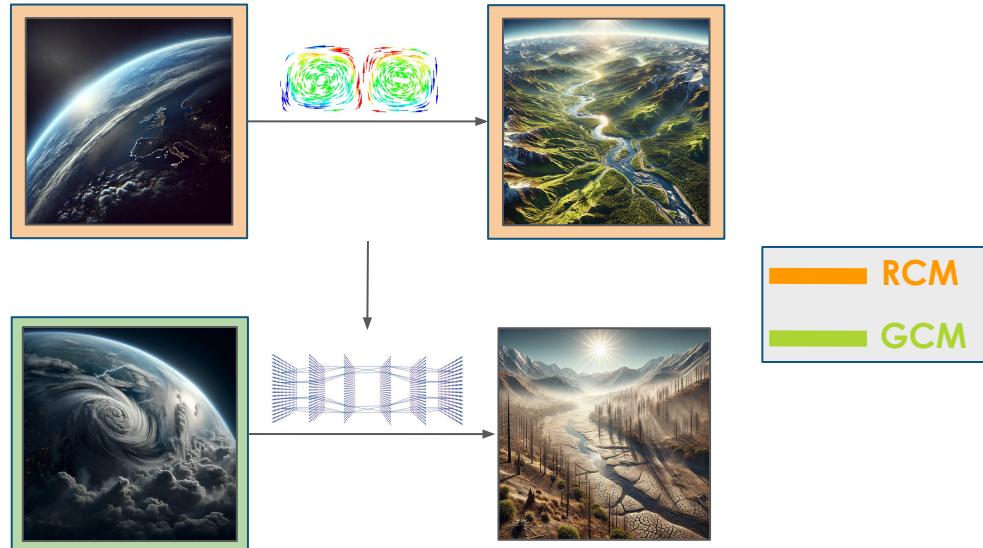
Emulated RCM  
projection

## Weather Turing Test

What are we actually learning from this approach to RCM emulation?

Large-scale  
variables

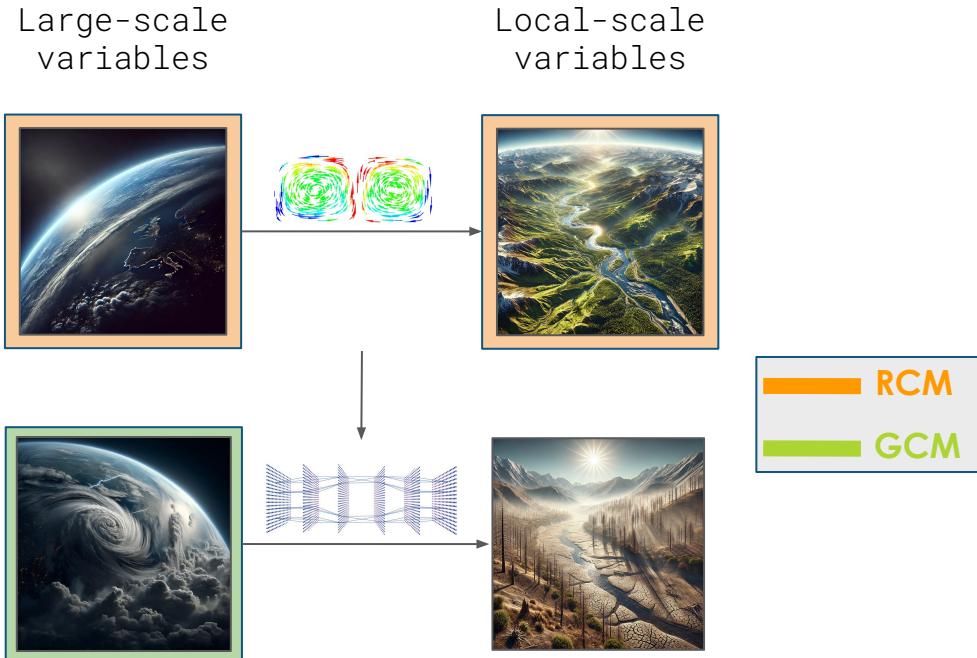
Local-scale  
variables



We are learning how the RCM **connects**  
**large-scale processes with local-scale**  
**dynamics.**

## Weather Turing Test

What are we actually learning from this approach to RCM emulation?



We are learning how the RCM **connects large-scale processes with local-scale dynamics.**

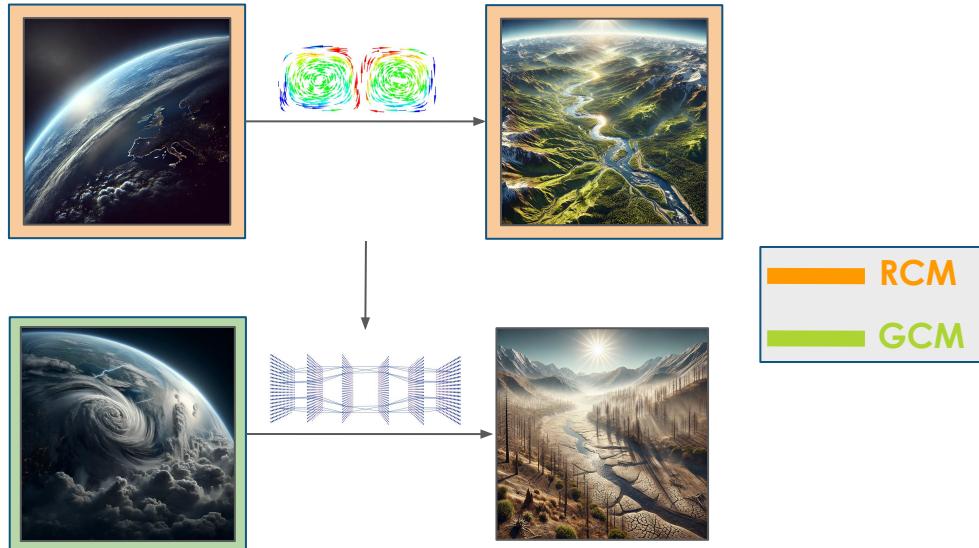
This allows us to obtain **high-resolution simulations for different GCMs** without running the RCM nested within each GCM.

## Weather Turing Test

What are we actually learning from this approach to RCM emulation?

Large-scale  
variables

Local-scale  
variables



We are learning how the RCM **connects large-scale processes with local-scale dynamics.**

Why do we use large-scale variables as predictor?

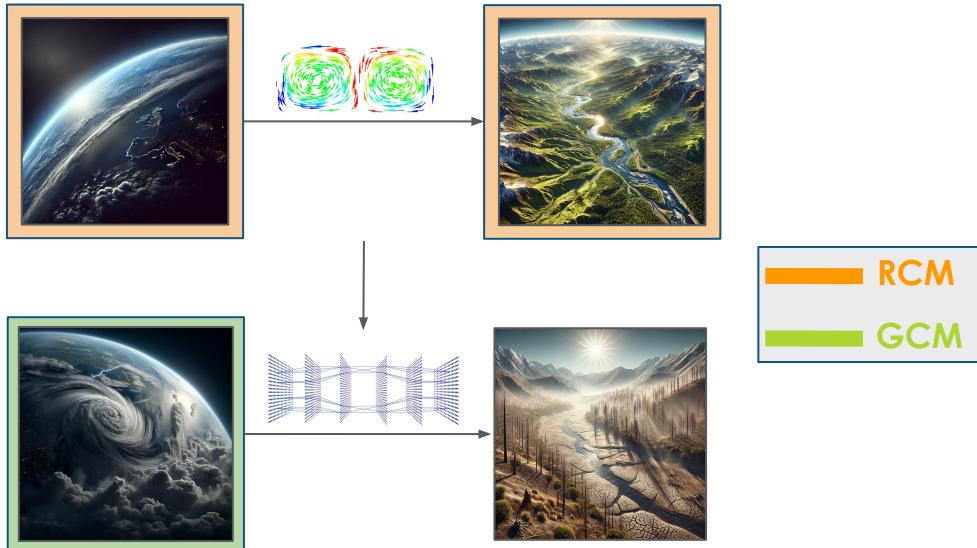
This allows us to obtain **high-resolution simulations for different GCMs** without running the RCM nested within each GCM.

## Weather Turing Test

What are we actually learning from this approach to RCM emulation?

Large-scale  
variables

Local-scale  
variables



We are learning how the RCM **connects large-scale processes with local-scale dynamics.**

**Why do we use large-scale variables as predictor?**

The typical coarse resolution of these variables is sufficient to capture their key physical phenomena.

This allows us to obtain **high-resolution simulations for different GCMs** without running the RCM nested within each GCM.

# Conclusions

- Global Climate Models (**GCMs**) are physical models that simulate the Earth's climate under different future scenarios.
  - **GCMs** work well at large scales but **struggle to capture local-scale detail.**
  - To address this, Regional Climate Models (**RCMs**) are **nested within GCMs** and run at **higher resolution over specific regions**. However, RCMs are computationally expensive.
- Deep Learning (**DL**) models can emulate the downscaling process (from coarse large-scale to fine local-scale), thus **emulating the RCM**.
  - Once trained, these emulators can be **applied to different GCM outputs**, enabling **fast and cost-efficient generation of high-resolution climate simulations.**

# Table of Contents

## **Part I: Review**

- DL for Weather/Climate
- RCM Emulation
- Conclusions

## **Part II: Training a RCM emulator**

- Problem Statement
- Training Frameworks
- Training and Inference

## **Part III: Evaluating a RCM emulator**

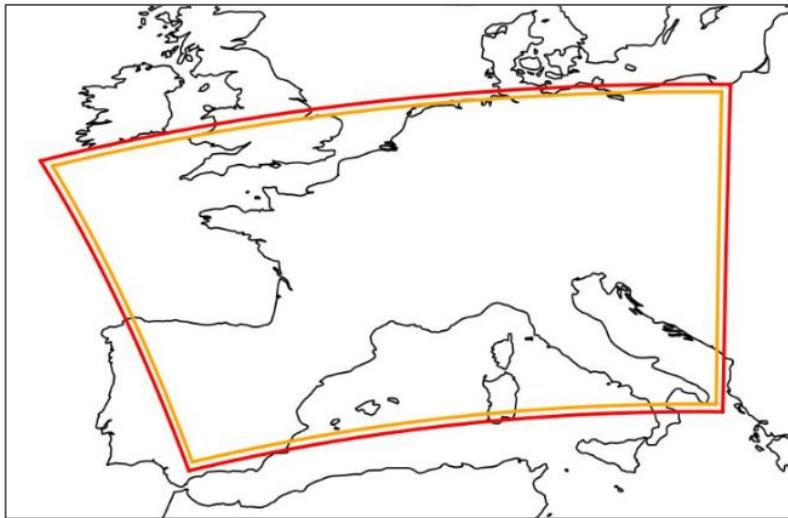
- Importance of Proper Evaluation
- Soft-Transferability
- Hard-Transferability
- Metrics

## **Part IV: The Importance of Benchmarks**

- What is a benchmark?
- CORDEX-ML-Bench

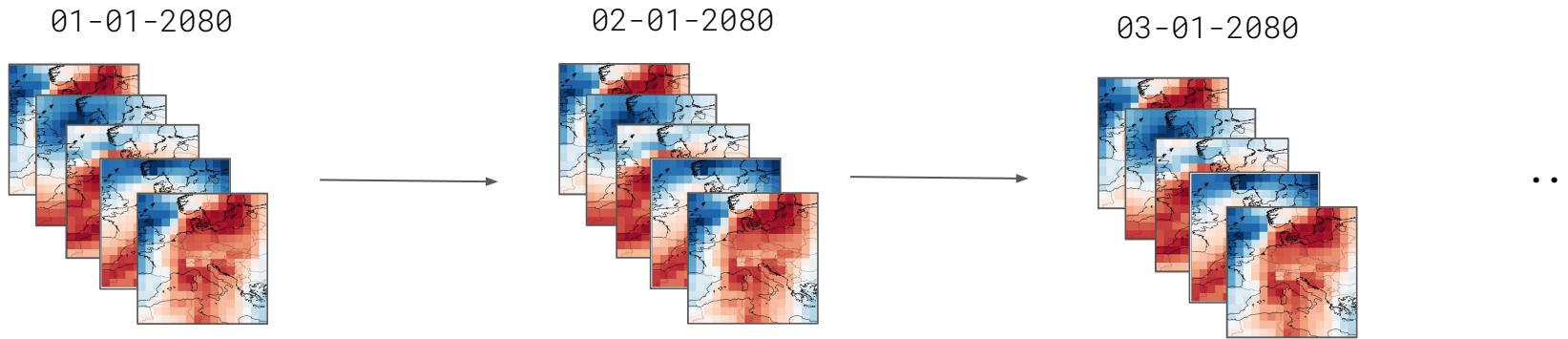
## Region of Interest

Suppose we need to study the **effects of climate change on future scenarios** for the **Alps region**.



## GCM Simulations

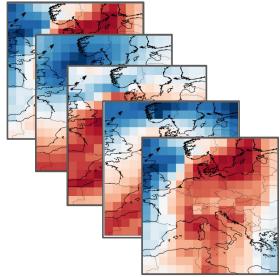
To do this, we obtain simulations of **future climate change scenarios from a GCM**, which we will denote as **GCM\_1**.



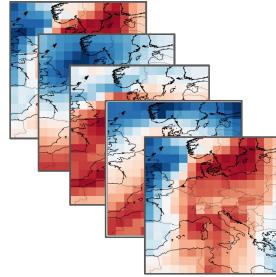
## GCM Simulations

To do this, we obtain simulations of **future climate change scenarios from a GCM**, which we will denote as **GCM\_1**.

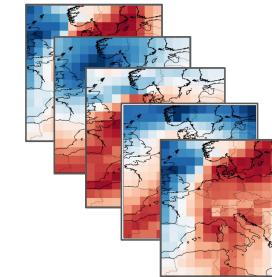
01-01-2080



02-01-2080

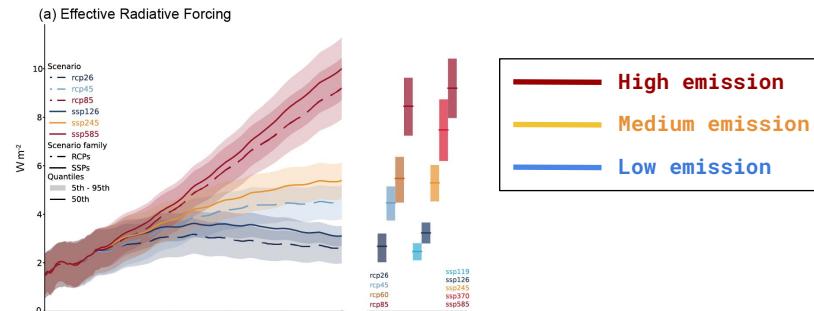


03-01-2080



...

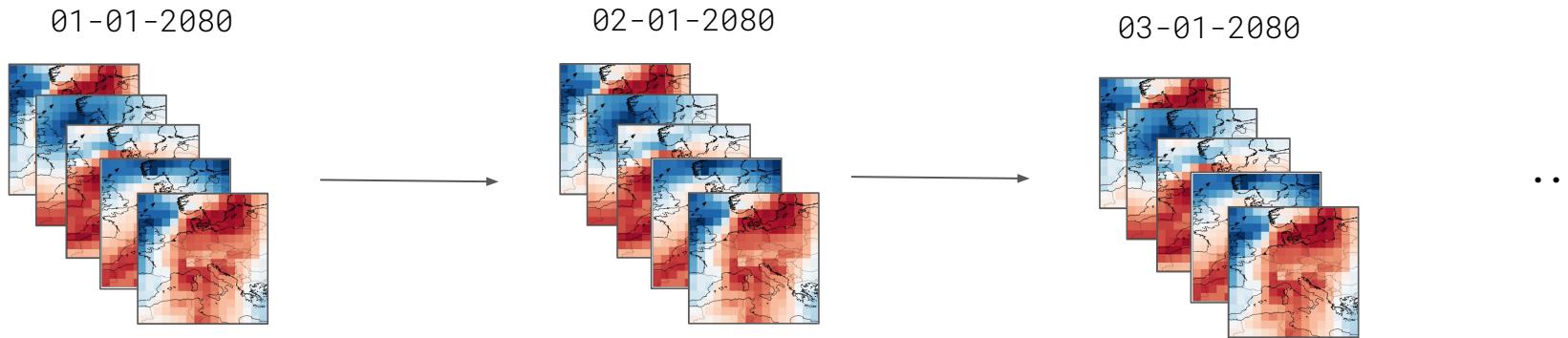
When running future simulations, we assume specific **emissions scenarios**



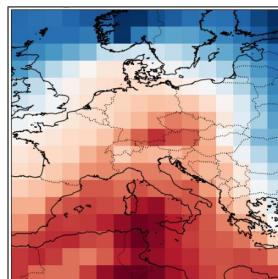
Source: IPCC WGI AR6, Ch4.

## GCM Simulations

To do this, we obtain simulations of **future climate change scenarios from a GCM**, which we will denote as **GCM\_1**.

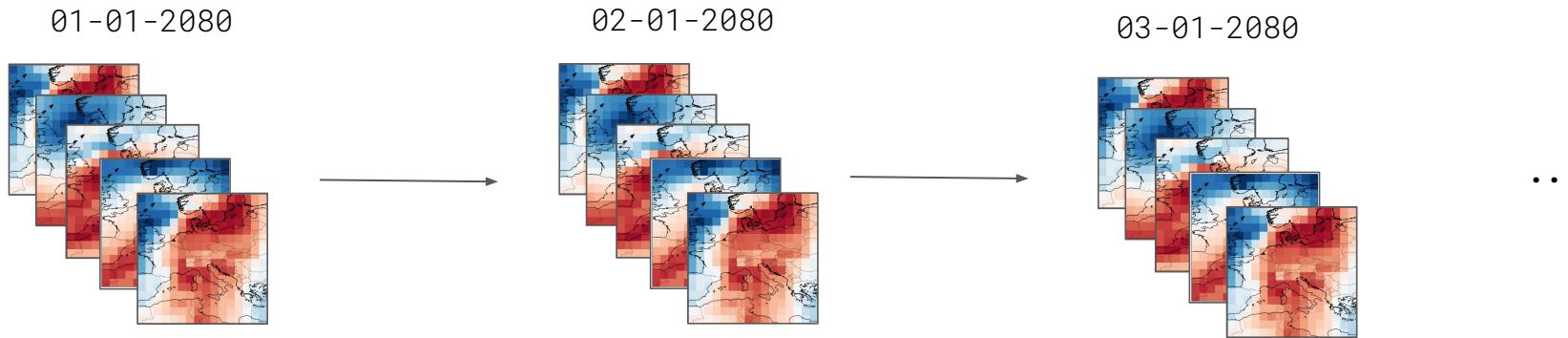


Unfortunately, due to the limitations of GCMs, the **simulations from GCM\_1 cannot accurately reproduce local-scale conditions**

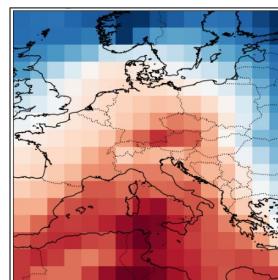


## GCM Simulations

To do this, we obtain simulations of **future climate change scenarios from a GCM**, which we will denote as **GCM\_1**.



Unfortunately, due to the limitations of GCMs, the **simulations from GCM\_1 cannot accurately reproduce local-scale conditions**

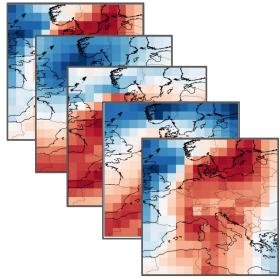


We need to run a **RCM**!

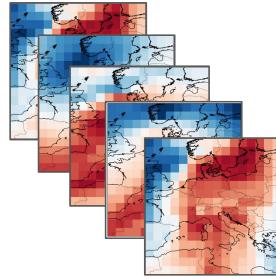
## RCM Simulations

We run an RCM using the boundary conditions from GCM\_1, thereby obtaining local-scale simulations of the future climate for the Alps region.

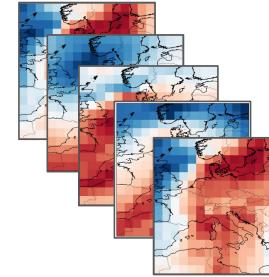
01-01-2080



02-01-2080



03-01-2080

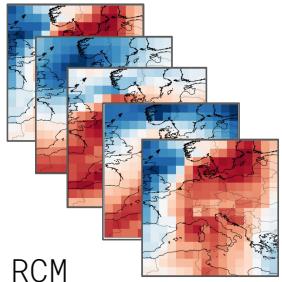


...

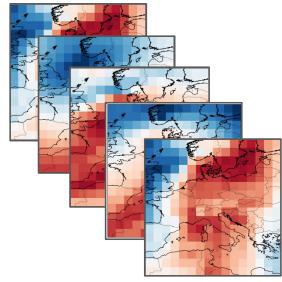
## RCM Simulations

We run an RCM using the boundary conditions from GCM\_1, thereby obtaining local-scale simulations of the future climate for the Alps region.

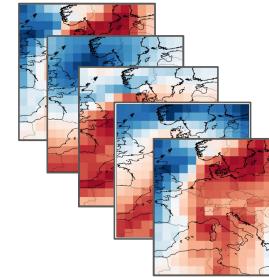
01-01-2080



02-01-2080

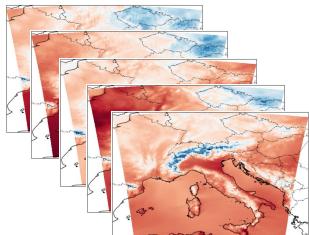
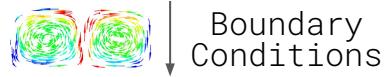


03-01-2080



...

RCM

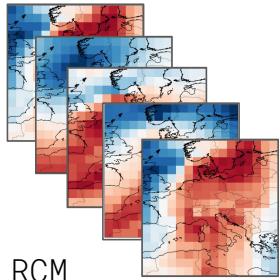


01-01-2080

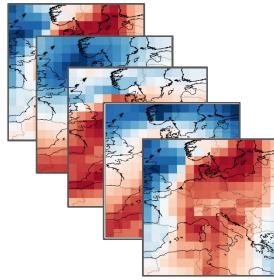
## RCM Simulations

We run an RCM using the boundary conditions from GCM\_1, thereby obtaining local-scale simulations of the future climate for the Alps region.

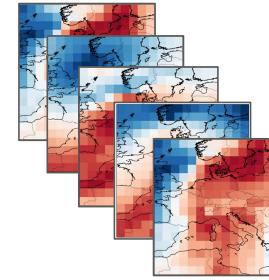
01-01-2080



02-01-2080



03-01-2080

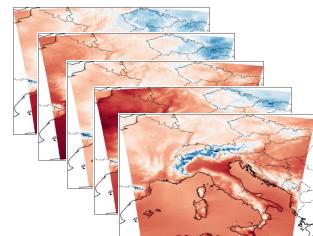


...

RCM

Boundary Conditions

Boundary Conditions



RCM



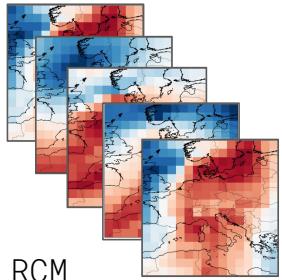
01-01-2080

01-01-2080

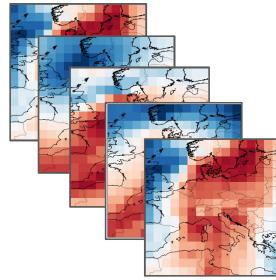
## RCM Simulations

We run an RCM using the boundary conditions from GCM\_1, thereby obtaining local-scale simulations of the future climate for the Alps region.

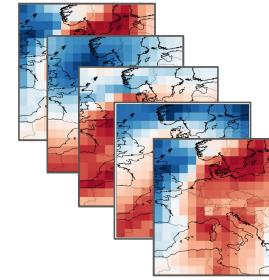
01-01-2080



02-01-2080



03-01-2080



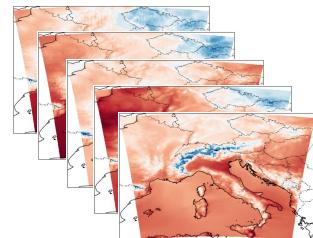
...

RCM

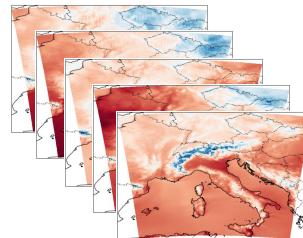
Boundary Conditions

Boundary Conditions

Boundary Conditions



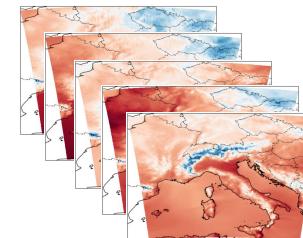
RCM



01-01-2080



RCM

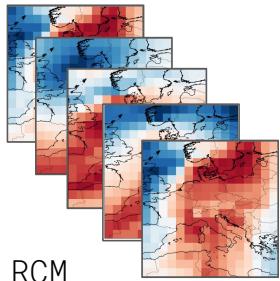


01-01-2080

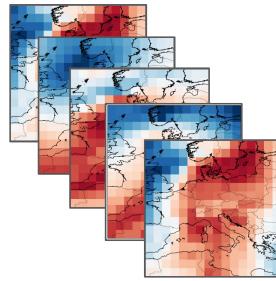
## RCM Simulations

We run an RCM using the boundary conditions from GCM\_1, thereby obtaining local-scale simulations of the future climate for the Alps region.

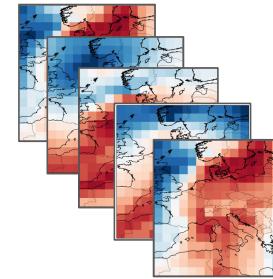
01-01-2080



02-01-2080



03-01-2080



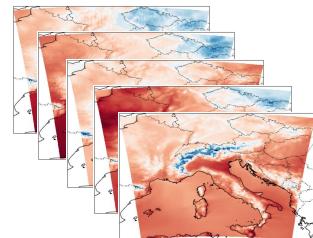
...

RCM

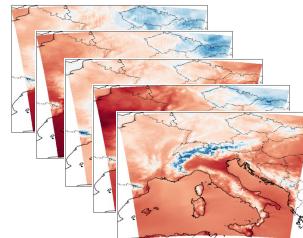
Boundary Conditions

Boundary Conditions

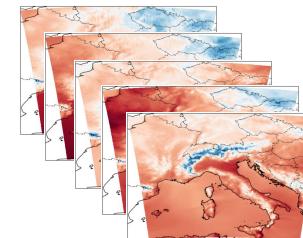
Boundary Conditions



RCM



RCM



RCM

01-01-2080



01-01-2080

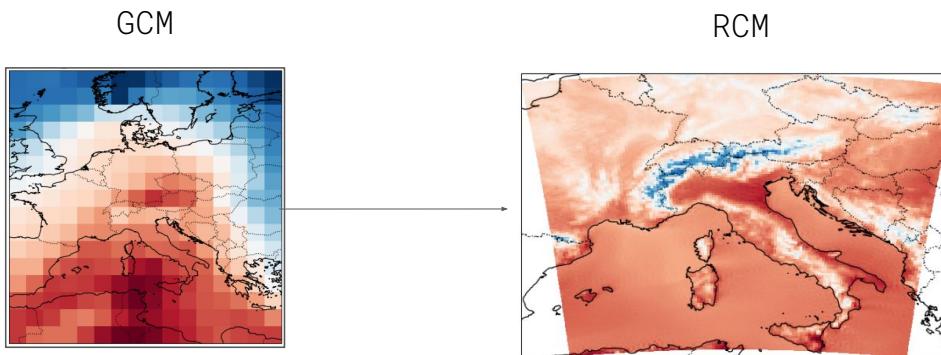
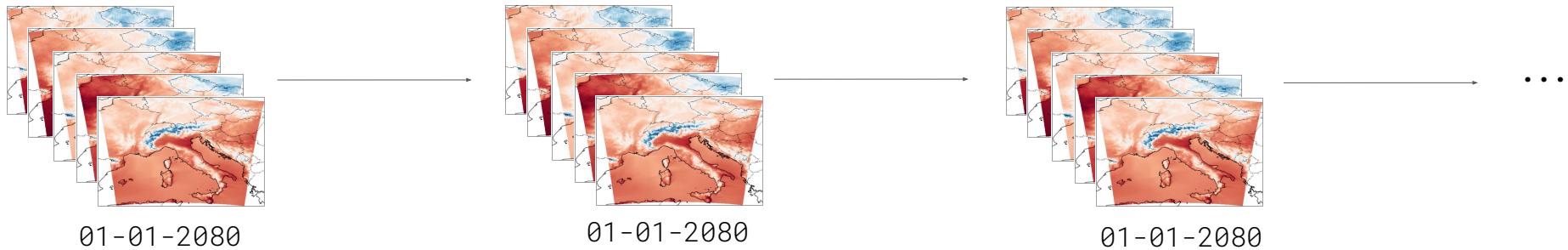


01-01-2080



## RCM Simulations

After **running the RCM** (with **GCM\_1** as **boundary conditions**), we obtained local-scale simulations for the Alps region.



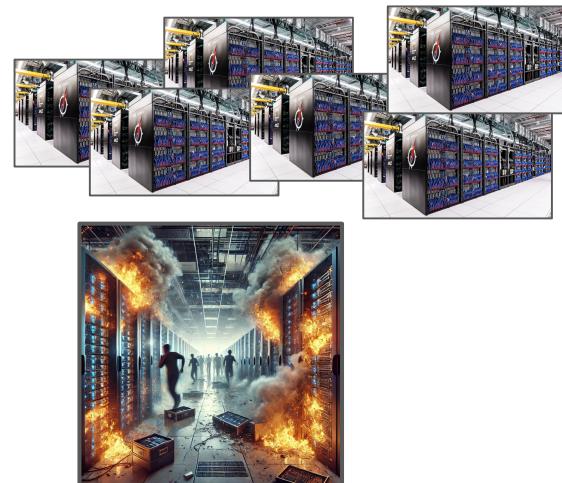
## RCM Emulation

But what if we want to run the RCM with **other GCMs** or under **different future scenarios?**

We need to re-run the RCM with  
different GCMs as boundary  
conditions



Computationally demanding



## RCM Emulation

But what if we want to run the RCM with **other GCMs** or under **different future scenarios**?

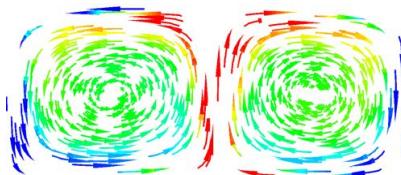
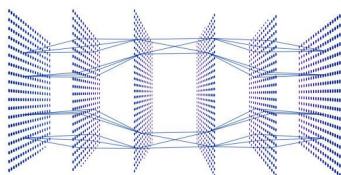
We need to re-run the RCM with  
different GCMs as boundary  
conditions



Computationally demanding



RCM Emulation



Deep Learning

RCM

## Perfect and Imperfect Frameworks

Our objective is to **emulate the downscaling function of the RCM** (run using GCM\_1 as boundary conditions) so that **it can be applied to different GCMs\_X**, thereby generating the corresponding simulations.

### Perfect Framework

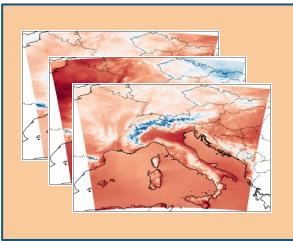
### Imperfect Framework



## Perfect and Imperfect Frameworks

Our objective is to **emulate the downscaling function of the RCM** (run using GCM\_1 as boundary conditions) so that **it can be applied to different GCMs\_X**, thereby generating the corresponding simulations.

### Perfect Framework



---

Large-scale  
variables

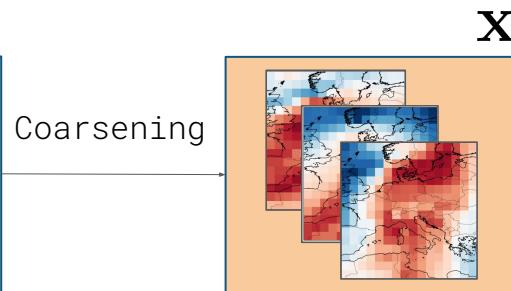
### Imperfect Framework



## Perfect and Imperfect Frameworks

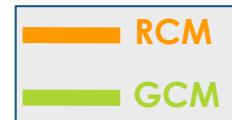
Our objective is to **emulate the downscaling function of the RCM** (run using GCM\_1 as boundary conditions) so that **it can be applied to different GCMs\_X**, thereby generating the corresponding simulations.

Perfect Framework



Large-scale  
variables

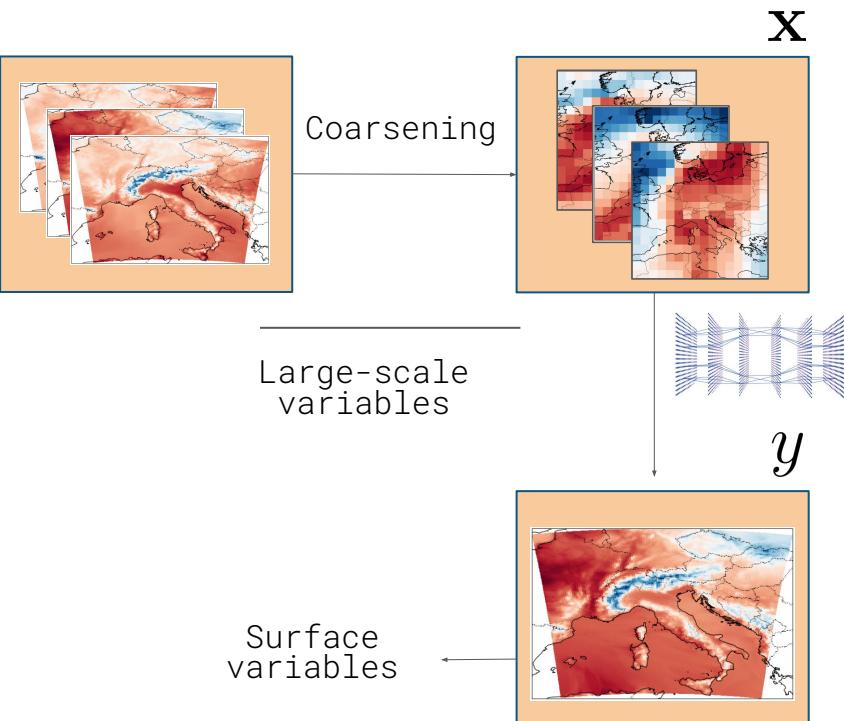
Imperfect Framework



## Perfect and Imperfect Frameworks

Our objective is to **emulate the downscaling function of the RCM** (run using GCM\_1 as boundary conditions) so that **it can be applied to different GCMs\_X**, thereby generating the corresponding simulations.

Perfect Framework



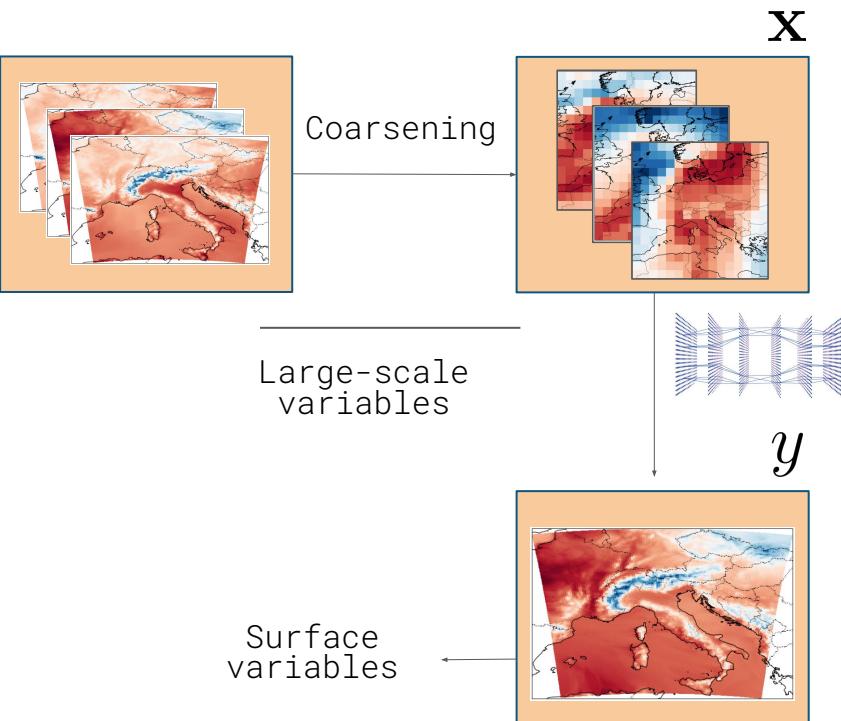
Imperfect Framework



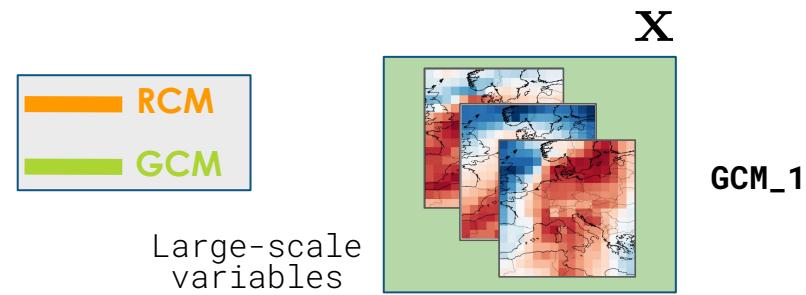
## Perfect and Imperfect Frameworks

Our objective is to **emulate the downscaling function of the RCM** (run using GCM\_1 as boundary conditions) so that **it can be applied to different GCMs\_X**, thereby generating the corresponding simulations.

Perfect Framework



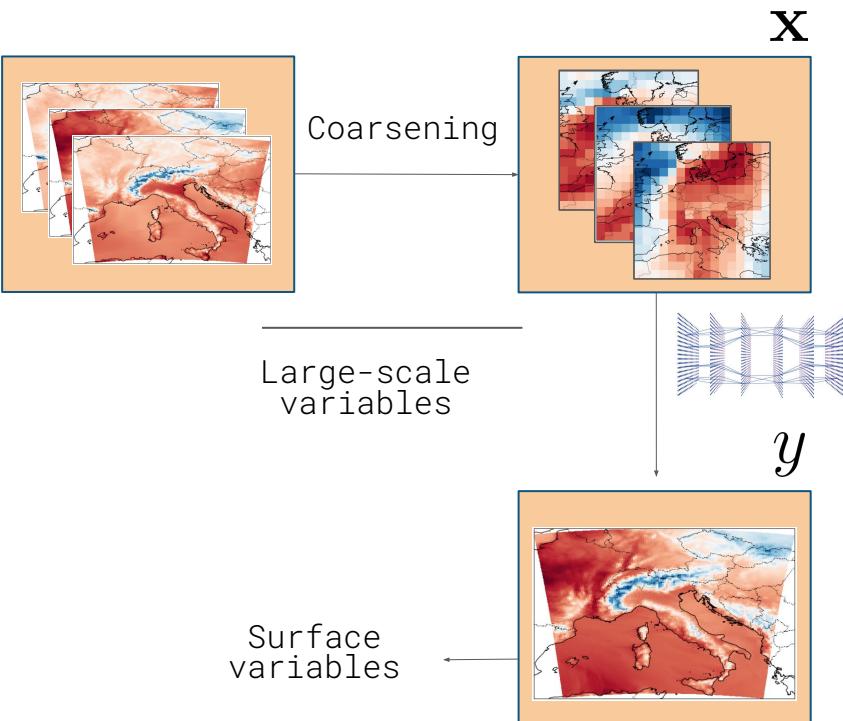
Imperfect Framework



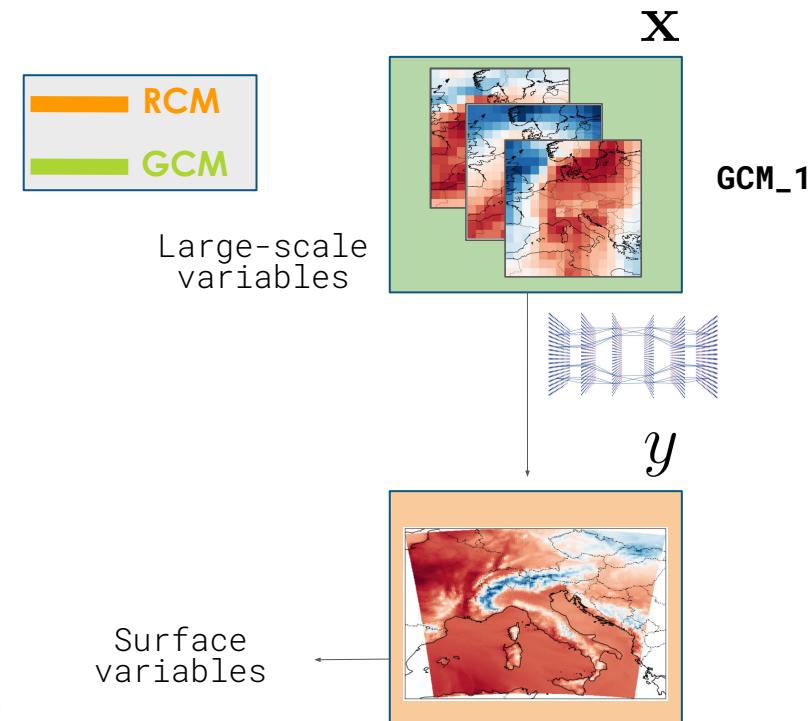
## Perfect and Imperfect Frameworks

Our objective is to **emulate the downscaling function of the RCM** (run using GCM\_1 as boundary conditions) so that **it can be applied to different GCMs\_X**, thereby generating the corresponding simulations.

**Perfect Framework**



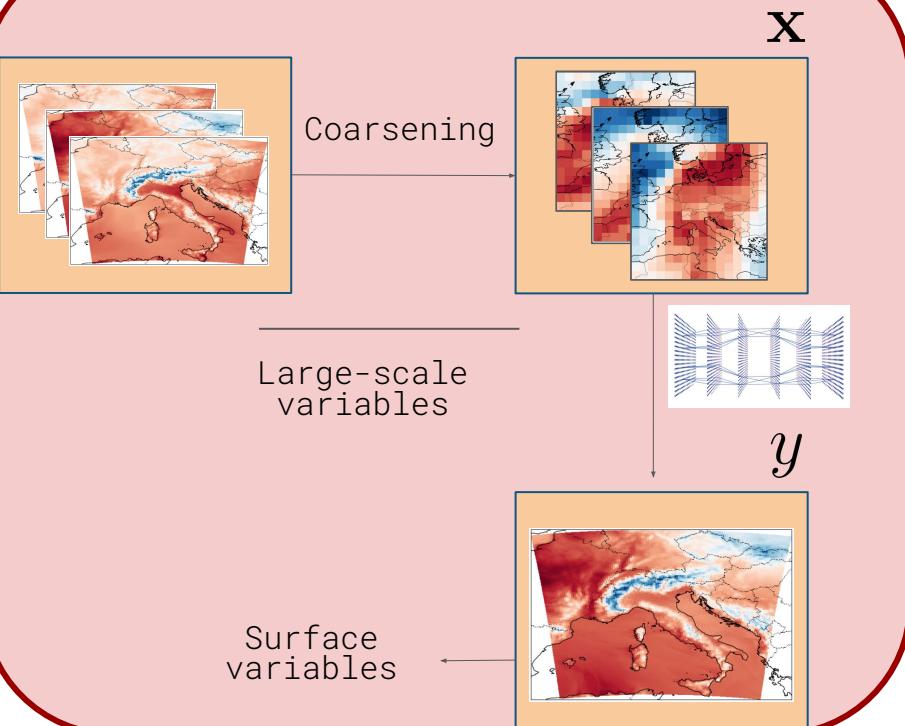
**Imperfect Framework**



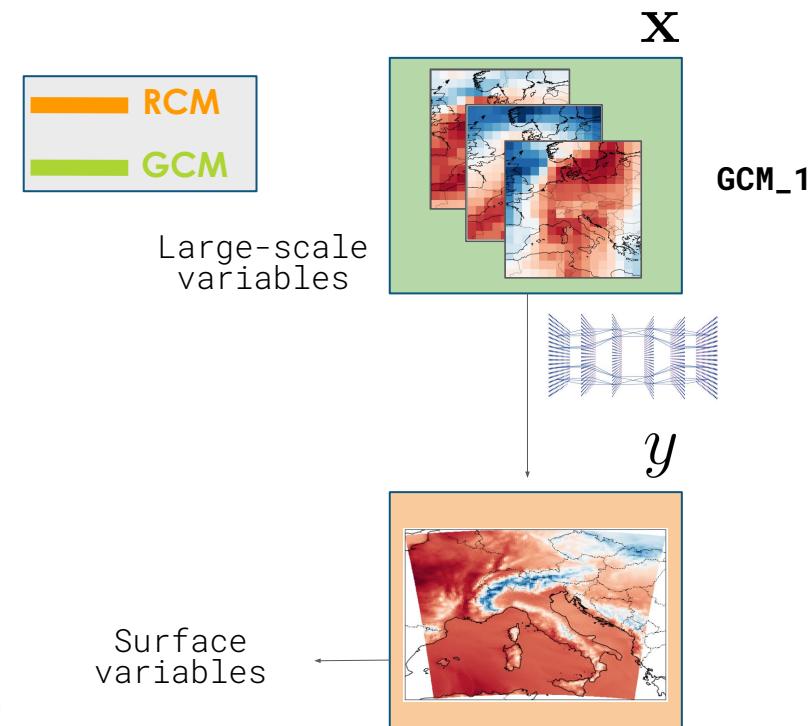
## Perfect and Imperfect Frameworks

Our objective is to **emulate the downscaling function of the RCM** (run using GCM\_1 as boundary conditions) so that **it can be applied to different GCMs\_X**, thereby generating the corresponding simulations.

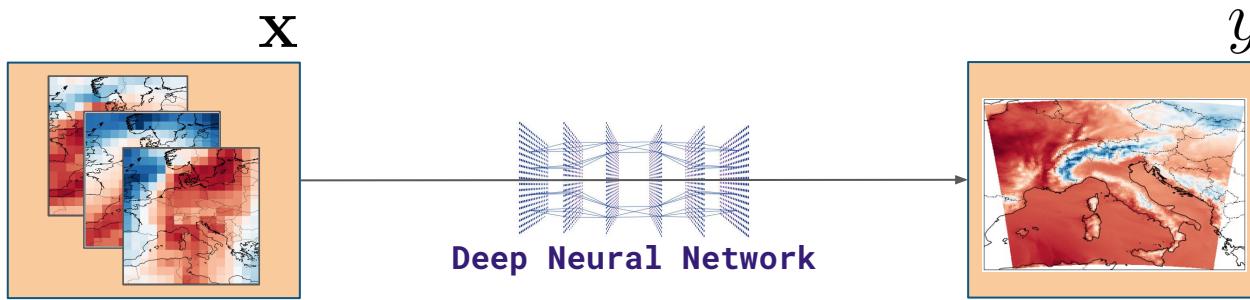
**Perfect Framework**



**Imperfect Framework**

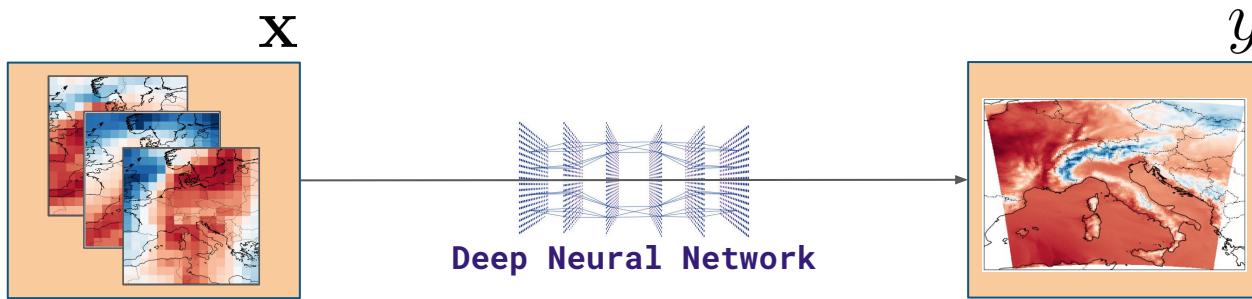


# Training



**Full dataset:** (1960-1980) and (2080-2100)

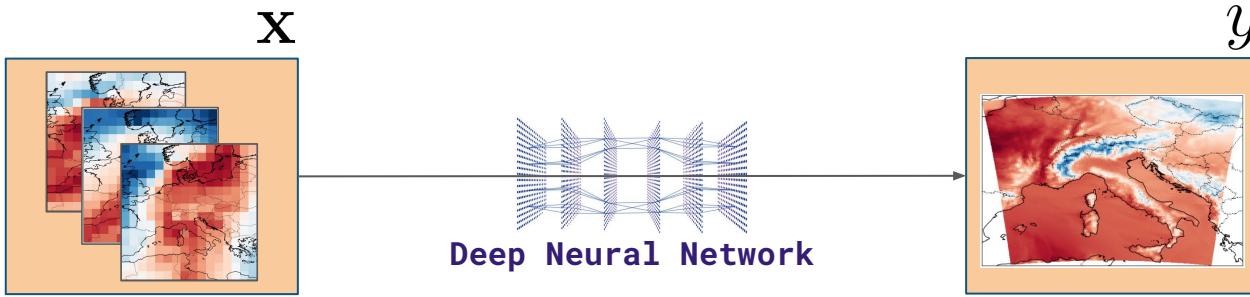
# Training



**Full dataset:** (1960-1980) and (2080-2100)

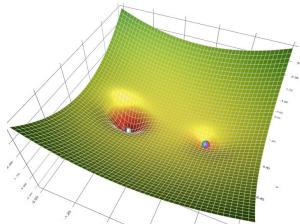


# Training

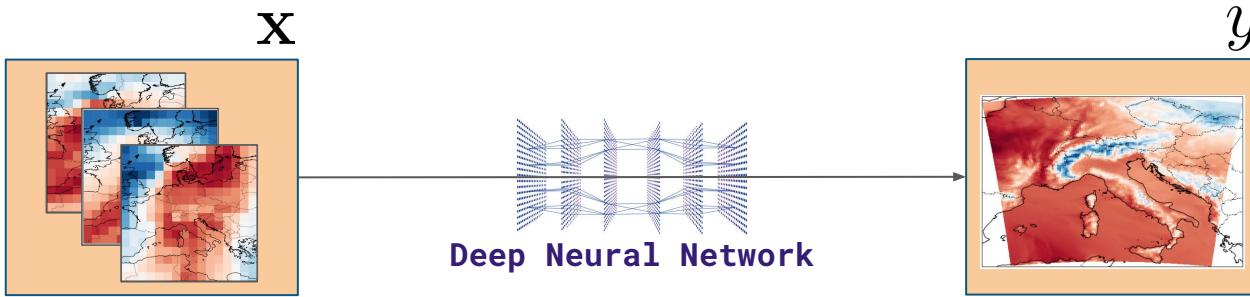


**Full dataset:** (1960-1980) and (2080-2100)

**Training set:** (1960-1975) and (2080-2095)



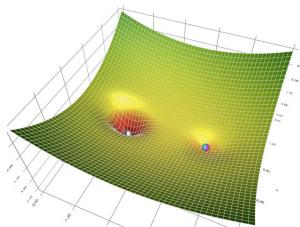
## Training



**Full dataset:** (1960-1980) and (2080-2100)

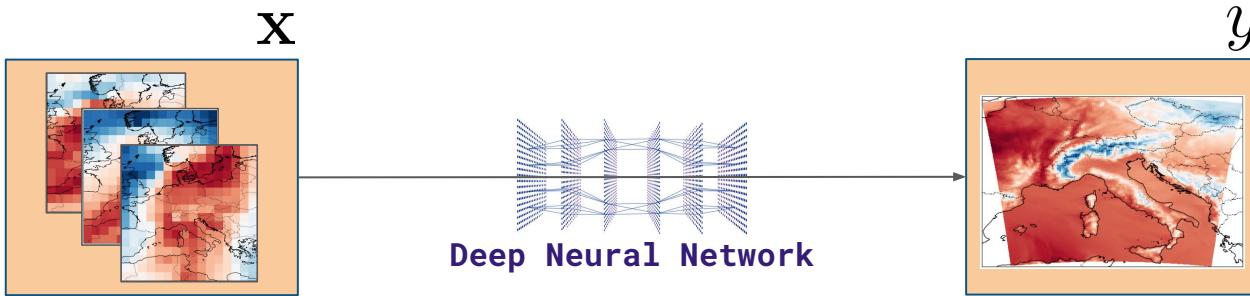
**Training set:** (1960-1975) and (2080-2095)

**Test set:** (1976-1980) and (2096-2100)

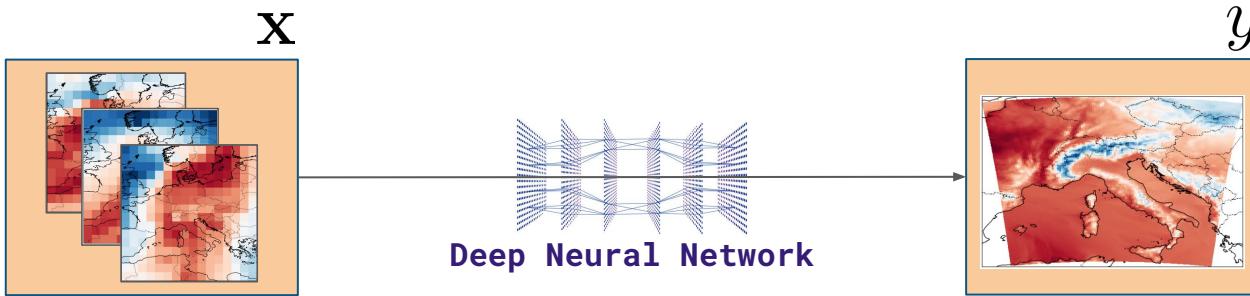


Is the model performing well on  
**unseen data?**

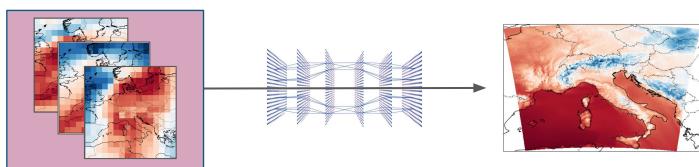
# Inference



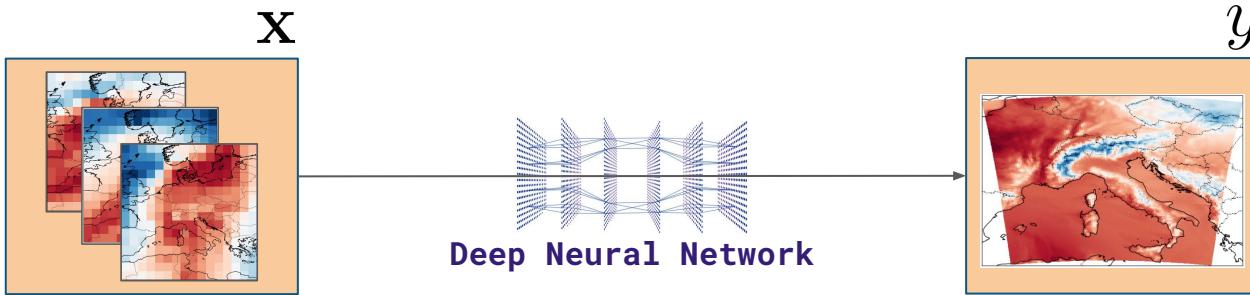
## Inference



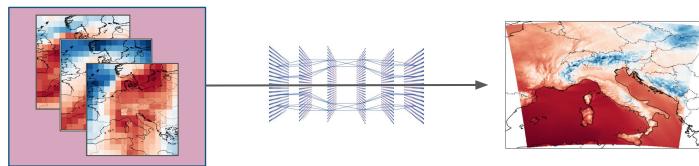
GCM\_2



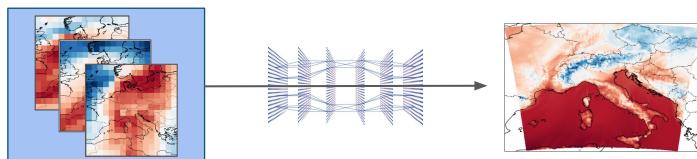
## Inference



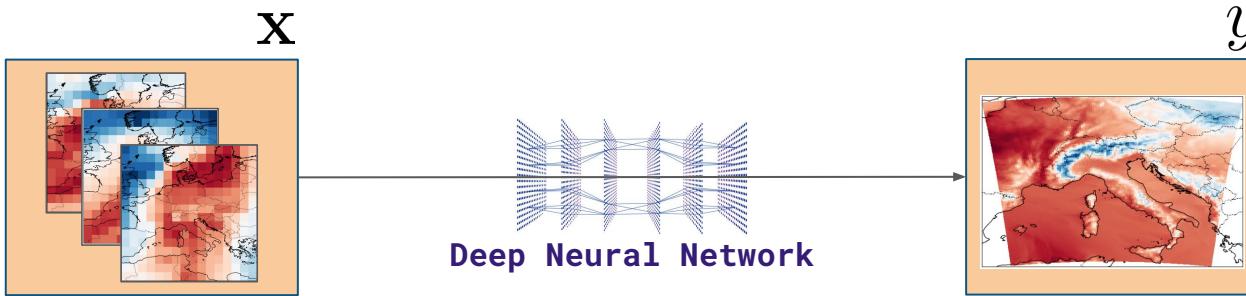
GCM\_2



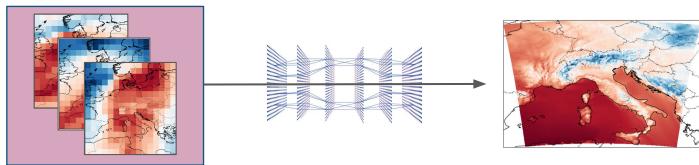
GCM\_3



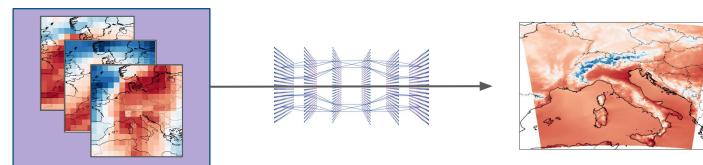
# Inference



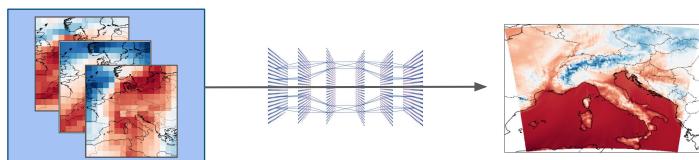
GCM\_2



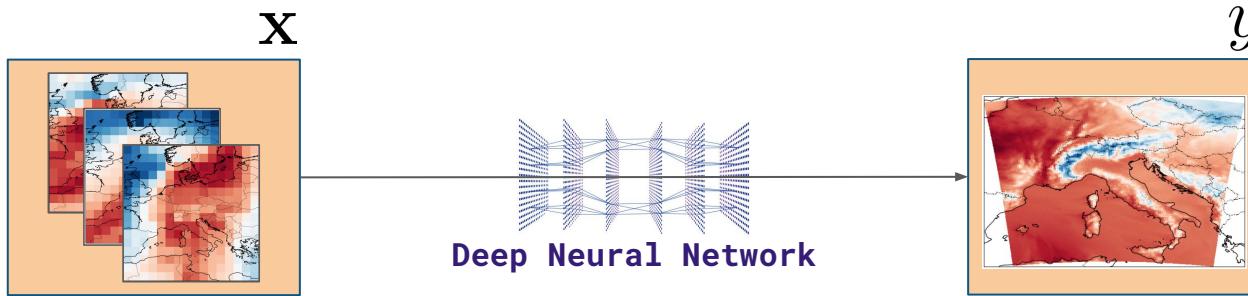
GCM\_4



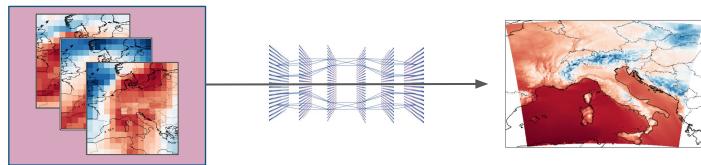
GCM\_3



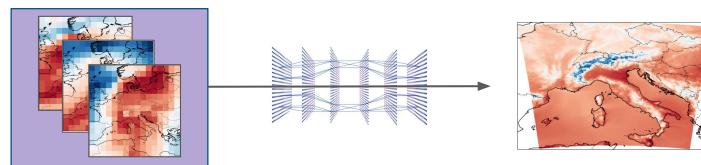
# Inference



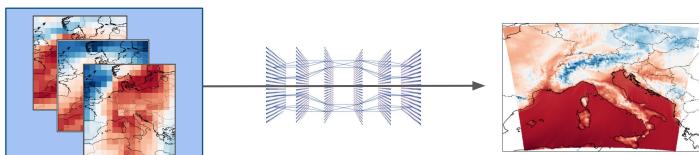
GCM\_2



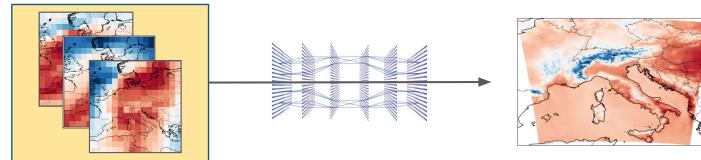
GCM\_4



GCM\_3



GCM\_5



# Table of Contents

## Part I: Review

- DL for Weather/Climate
- RCM Emulation
- Conclusions

## Part II: Training a RCM emulator

- Problem Statement
- Training Frameworks
- Training and Inference

## Part III: Evaluating a RCM emulator

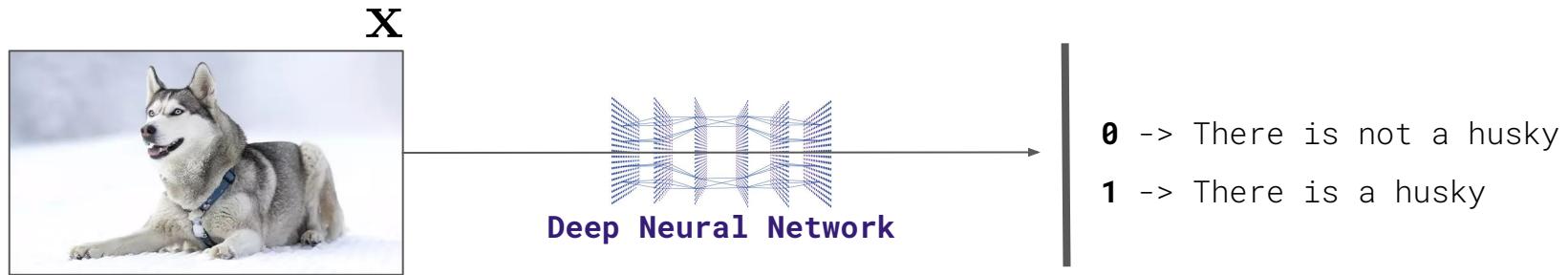
- Importance of Proper Evaluation
- Soft-Transferability
- Hard-Transferability
- Metrics

## Part IV: The Importance of Benchmarks

- What is a benchmark?
- CORDEX-ML-Bench

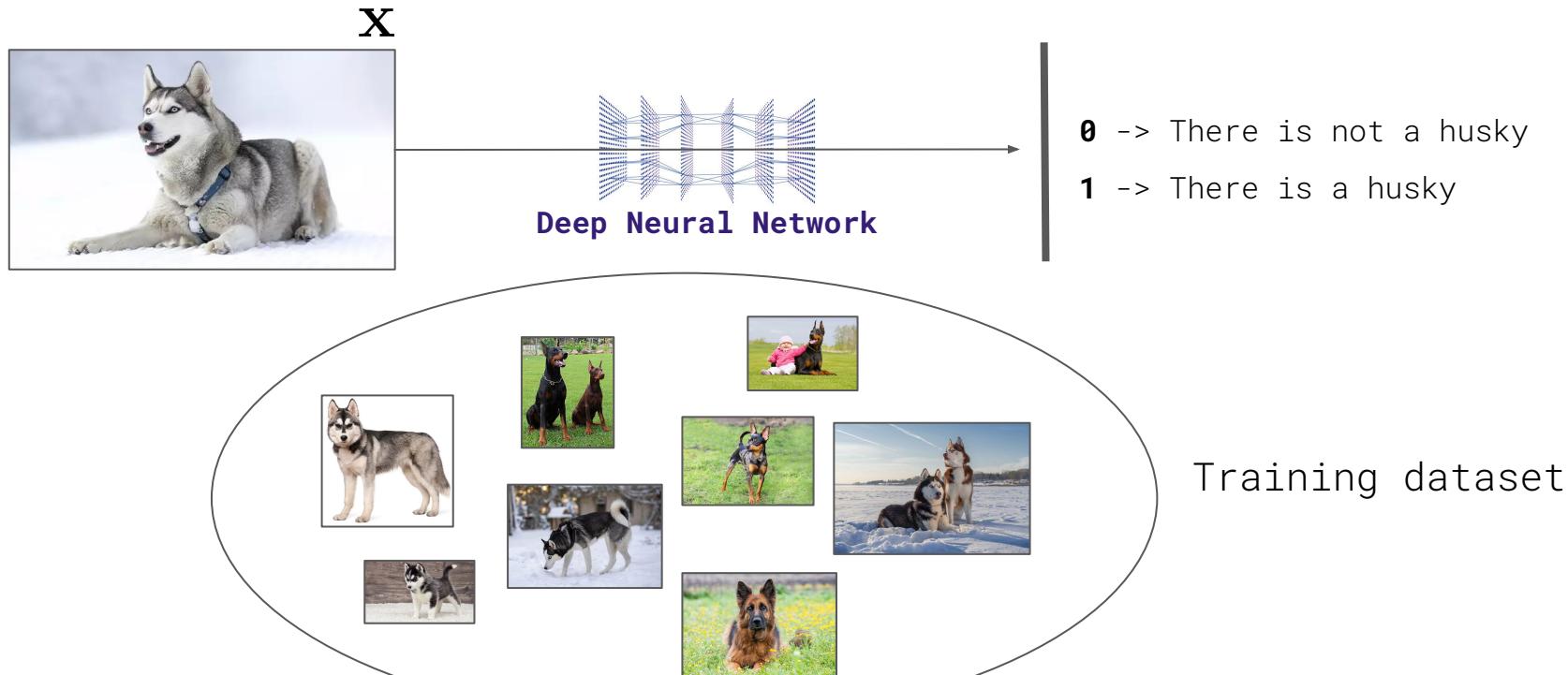
## Husky Classification

Suppose we want to train an **image classification model** to detect **Huskies** in images



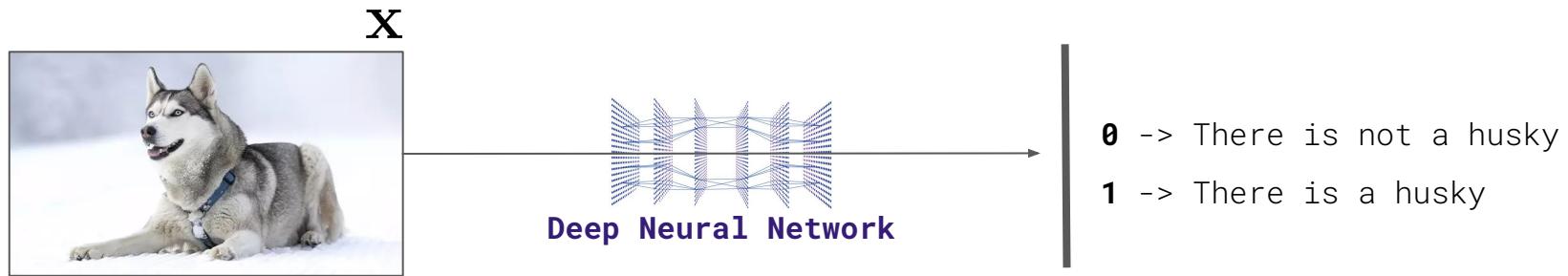
## Husky Classification

Suppose we want to train an **image classification model** to detect **Huskies** in images



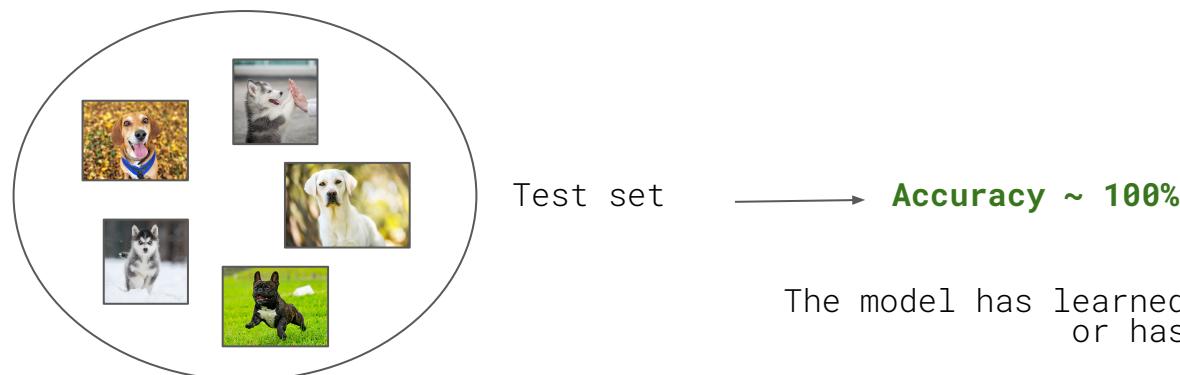
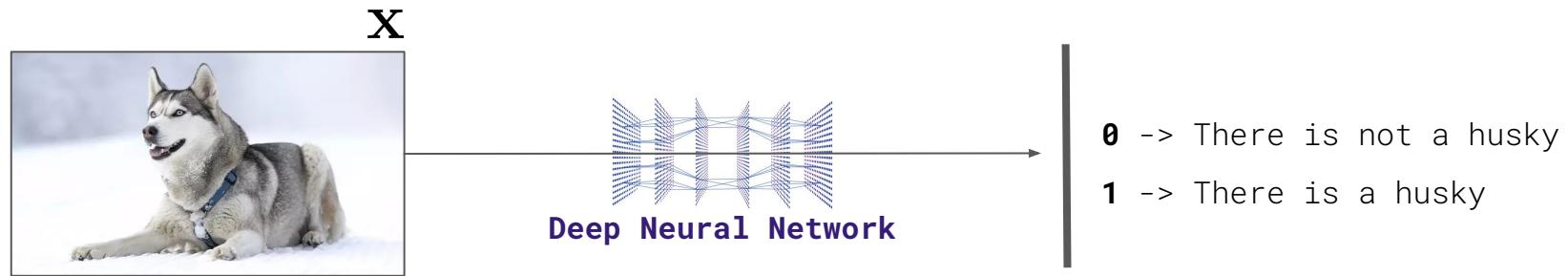
## Husky Classification

To assess its **accuracy**, we evaluate it on a **test set**



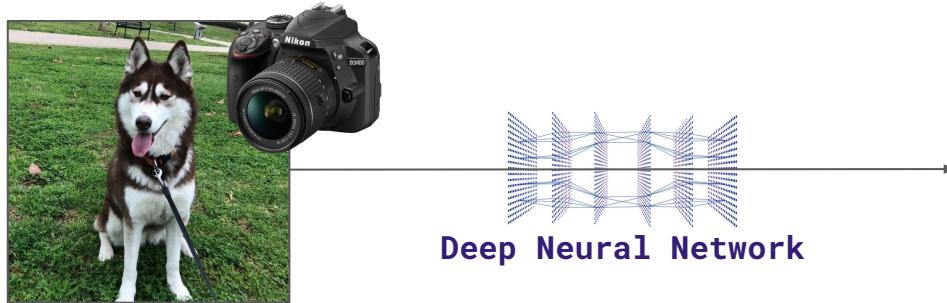
## Husky Classification

To assess its **accuracy**, we evaluate it on a **test set**



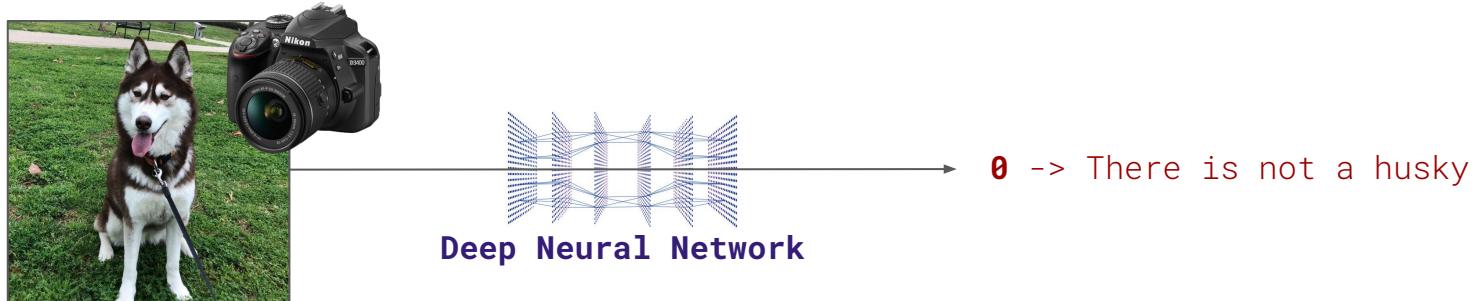
## Husky Classification

However, when we test the model on **real-world data**, it **doesn't seem to work well**



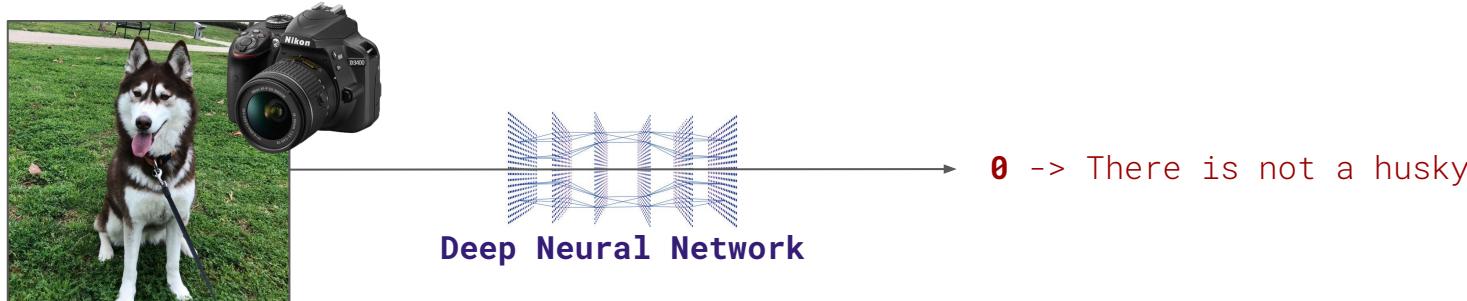
## Husky Classification

However, when we test the model on **real-world data**, it **doesn't seem to work well**



## Husky Classification

However, when we test the model on **real-world data**, it **doesn't seem to work well**



What is going on? Why is the model failing?



## Husky Classification

The model is not **learning** the concept of Husky, but the **concept of snow**



## Husky Classification

The model is not **learning** the concept of Husky, but the **concept of snow**



**Evaluation goes beyond using a partition of the training distribution**

## Husky Classification

The model is not **learning** the concept of Husky, but the **concept of snow**



**Evaluation goes beyond using a partition of the training distribution**

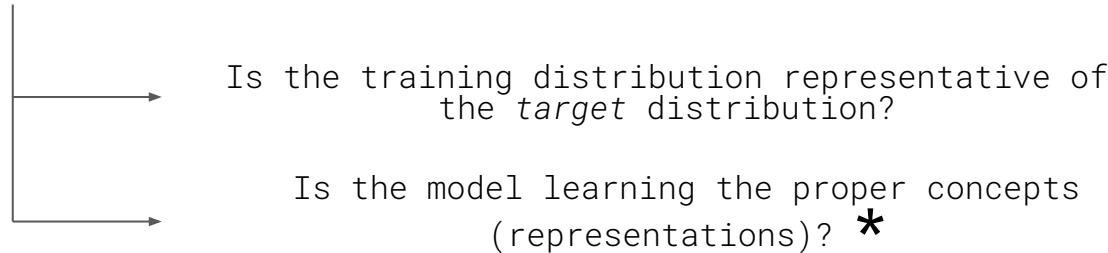
→ Is the training distribution representative of  
the *target* distribution?

## Husky Classification

The model is not **learning** the concept of Husky, but the **concept of snow**



**Evaluation goes beyond using a partition of the training distribution**



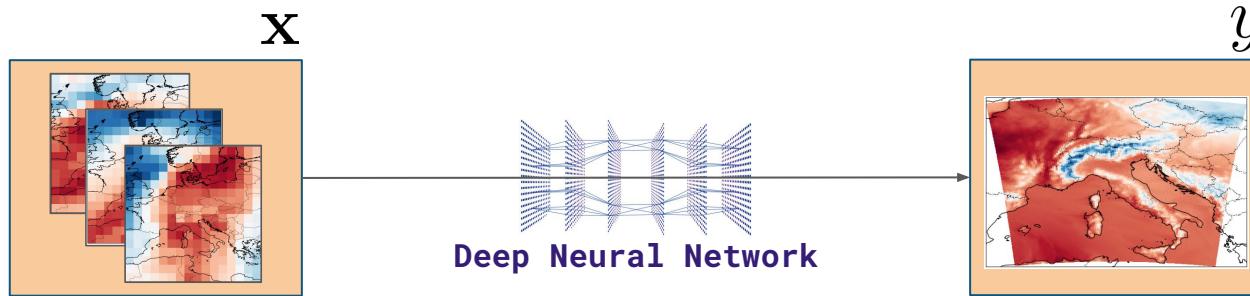
\* Or is it relying on shortcut learning?

## What about RCM emulation?

**Soft-transferability**

## What about RCM emulation?

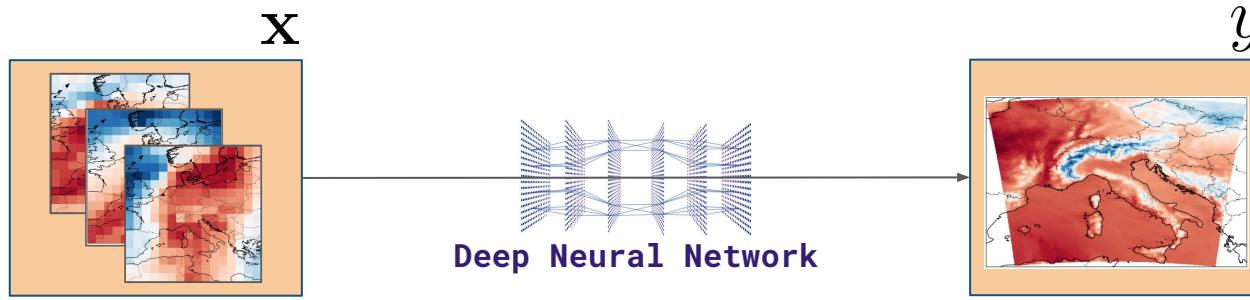
We have trained our DL model on a specific RCM ran with GCM\_1 as boundary conditions



**Full dataset:** (1960-1980) and (2080-2100)

## What about RCM emulation?

We have trained our DL model on a specific RCM ran with GCM\_1 as boundary conditions

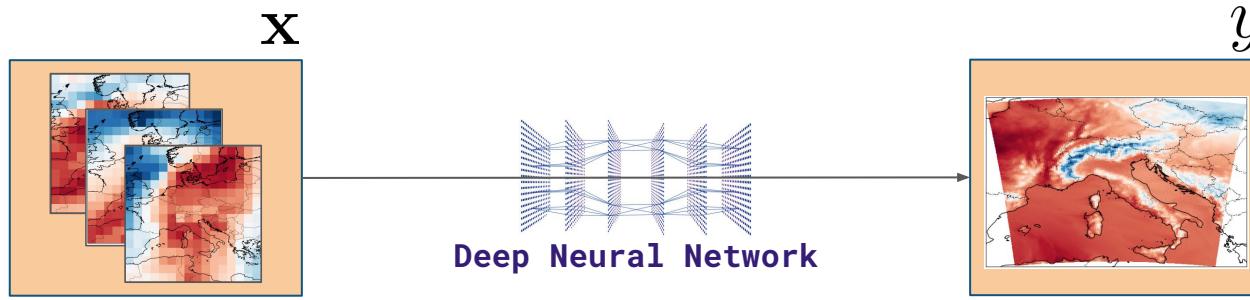


**Full dataset:** (1960-1980) and (2080-2100)

1. What if we simulate a **period not covered by the full dataset?**
2. What if we simulate an **emissions scenario different from that of GCM\_1** used to run the RCM?

## What about RCM emulation?

We have trained our DL model on a specific RCM ran with GCM\_1 as boundary conditions



**Full dataset:** (1960-1980) and (2080-2100)

1. What if we simulate a **period not covered by the full dataset?**
2. What if we simulate an **emissions scenario different from that of GCM\_1** used to run the RCM?

**Soft-transferability**

## Interpolation/Extrapolation

What if we simulate a **period not covered by the full dataset?**

This is not a major issue for **RCM emulation**, since **we can train the model using both historical and future data**, allowing us to just having to evaluate its performance for **interpolation**.

**Full dataset:** (1960-1980) and (2080-2100)

## Interpolation/Extrapolation

What if we simulate a **period not covered by the full dataset?**

This is not a major issue for **RCM emulation**, since **we can train the model using both historical and future data**, allowing us to just having to evaluate its performance for **interpolation**.

**Full dataset:** (1960-1980) and (2080-2100)

However, for other approaches like **observational statistical downscaling**, where **we have only historical data**, it is important to properly **evaluate how well they can extrapolate** to future scenarios.

OCTOBER 2025

GONZÁLEZ-ABAD AND GUTIÉRREZ

### Are Deep Learning Methods Suitable for Downscaling Global Climate Projections? An Intercomparison for Temperature and Precipitation over Spain

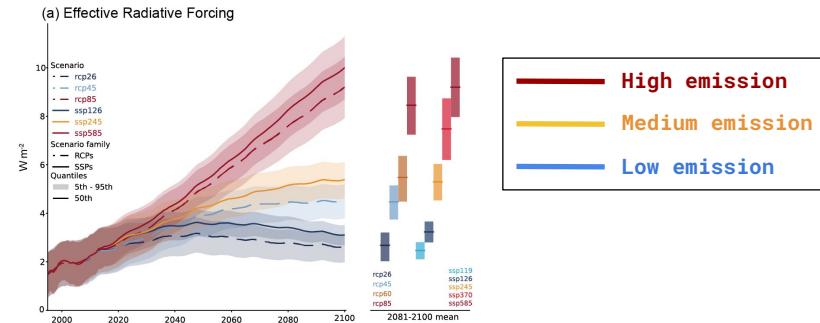
JOSE GONZÁLEZ-ABAD<sup>a</sup> AND JOSÉ MANUEL GUTIÉRREZ<sup>a</sup>

<sup>a</sup> *Instituto de Física de Cantabria (IFCA), CSIC-Universidad de Cantabria, Santander, Spain*

## Different Scenarios

What if we simulate an **emissions scenario different from that of GCM\_1** used to run the RCM?

Suppose the RCM simulations used for training were run assuming a **low emissions scenario**



What if I run the **inference on GCM\_1**, but this time it was run under a **high emissions scenario?**

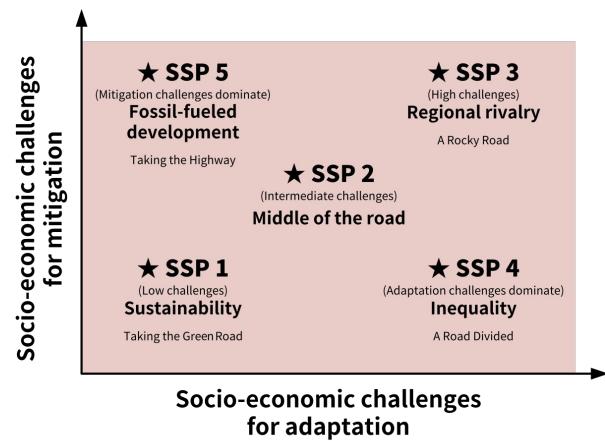


Is the RCM emulator able to extrapolate to climates not seen during training?

## Different Scenarios

What if we simulate an **emissions scenario different from that of GCM\_1** used to run the RCM?

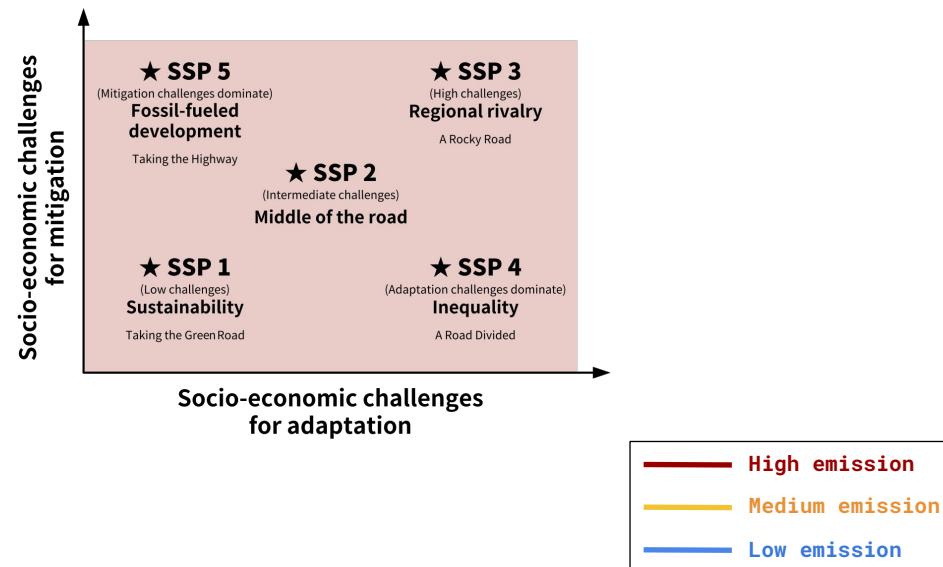
Assessing this property is key if we want to fill the GCM/RCM matrix across different scenarios



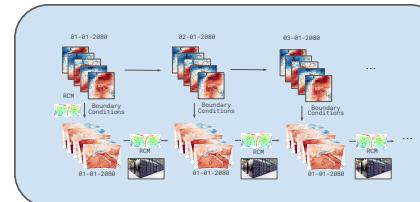
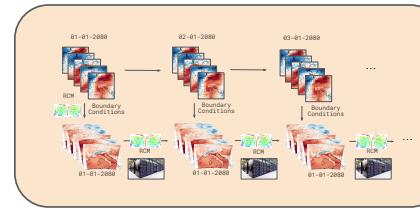
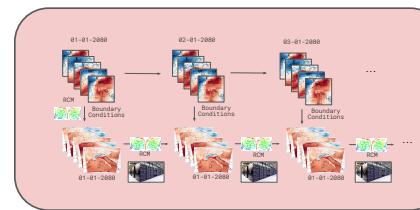
## Different Scenarios

What if we simulate an **emissions scenario different from that of GCM\_1** used to run the RCM?

Assessing this property is key if we want to fill the GCM/RCM matrix across different scenarios



For this evaluation, we need **RCM simulations** that use GCM\_1 as boundary conditions, **with different emissions scenarios**

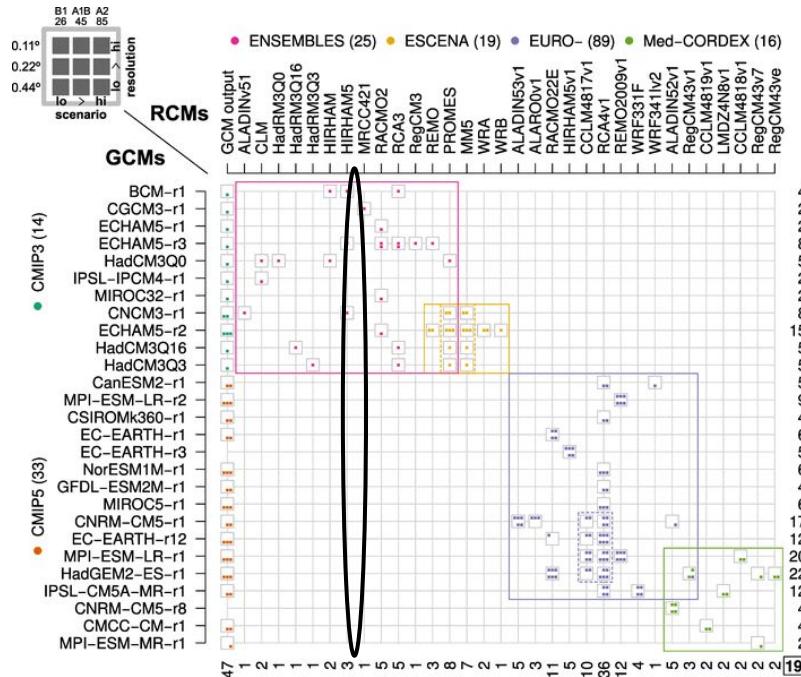


## Different GCMs

**Hard-transferability**

# Different GCMs

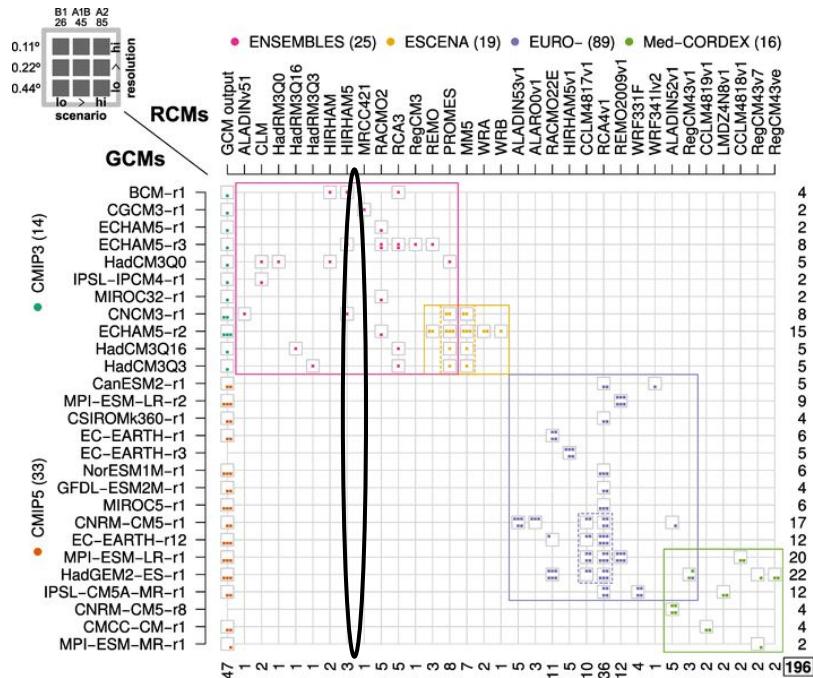
Our RCM emulator has been trained with simulations of an RCM run with GCM\_1; however, **the main purpose of an RCM emulator is to generate simulations for different GCMs**



Fuente: Fernández, J., Frías, M. D., Cabos, W. D., Cofiño, A. S., Domínguez, M., Fita, L., ... & Sánchez, E. (2019). Consistency of climate change projections from multiple global and regional model intercomparison projects. Climate dynamics, 52, 1139-1156.

# Different GCMs

Our RCM emulator has been trained with simulations of an RCM run with GCM\_1; however, **the main purpose of an RCM emulator is to generate simulations for different GCMs**



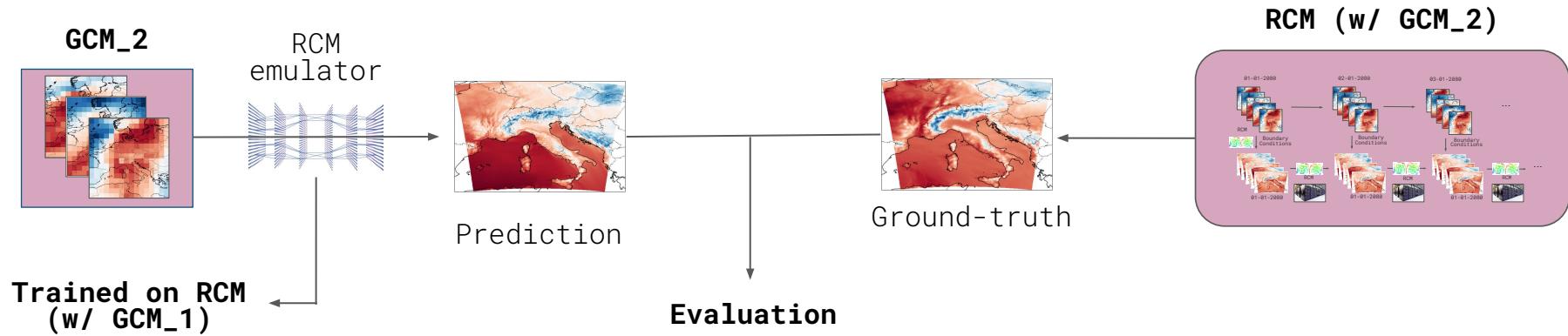
How can we evaluate the **performance** of the RCM emulator **across different GCMs**?



**Hard-transferability**

## Different GCMs

To evaluate the **hard-transferability** of a RCM emulator we need to run that same RCM but over a **different GCM** (for instance **GCM\_2**)



Is the RCM emulator able to extrapolate to other GCMs?

## Recap

### Soft-transferability

Is the RCM emulator able to interpolate/extrapolate in time?

Is the RCM emulator able to extrapolate to climates (emission scenarios) not seen during training?

### Hard-transferability

Is the RCM emulator able to extrapolate to other GCMs?

## Recap

### Soft-transferability

Is the RCM emulator able to interpolate/extrapolate in time?

Is the RCM emulator able to extrapolate to climates (emission scenarios) not seen during training?

### Hard-transferability

Is the RCM emulator able to extrapolate to other GCMs?

Unfortunately, most evaluations require **re-running the RCM we aim to emulate**, which is precisely what RCM emulation seeks to avoid



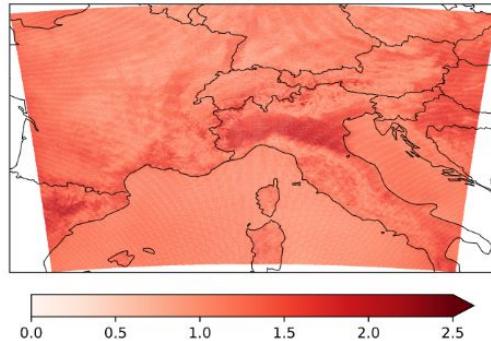
The proper evaluation of RCM emulator is challenging

## ML-based Metrics

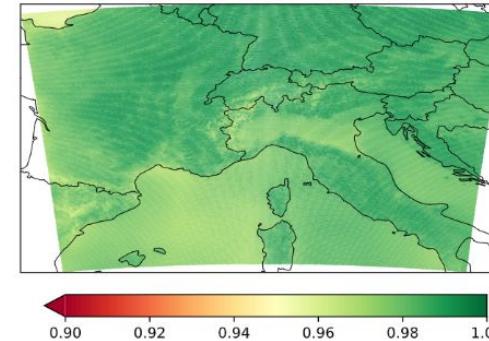
So far, we have reviewed the evaluation frameworks (i.e., what data to use for assessing different aspects of the RCM emulator), but **what metrics should we compute?**

**Classic ML metrics** (point-based)

**Root Mean Squared Error  
(RMSE)**



**Pearson Correlation**

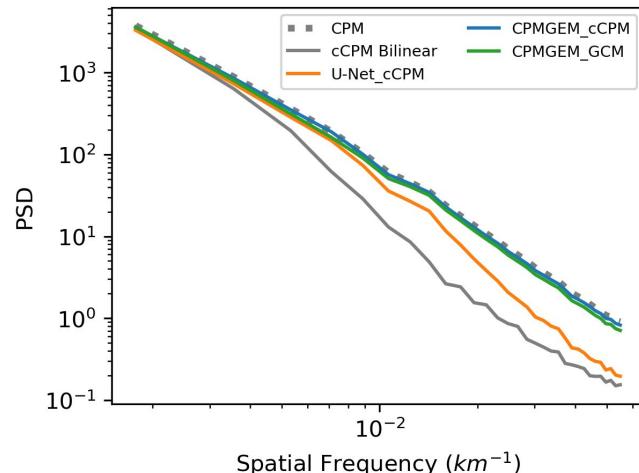


## ML-based Metrics

So far, we have reviewed the evaluation frameworks (i.e., what data to use for assessing different aspects of the RCM emulator), but **what metrics should we compute?**

### Classic ML metrics

#### Power Spectral Density (PSD)



Is the RCM emulator producing  
**smoothed/blurry** predictions?

Ground-truth: RCM  
Prediction: RCM emulator

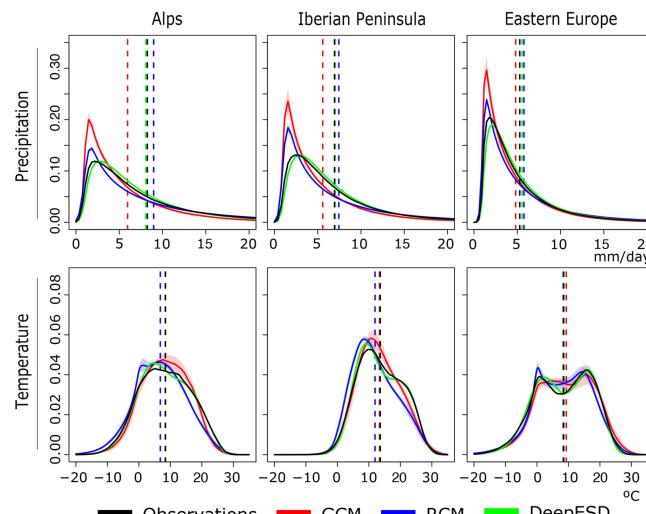
**Source:** Addison, H., Kendon, E., Ravuri, S., Aitchison, L., & Watson, P. A. (2024). Machine learning emulation of precipitation from km-scale regional climate simulations using a diffusion model. arXiv preprint arXiv:2407.14158.

## Distributional-based Metrics

So far, we have reviewed the evaluation frameworks (i.e., what data to use for assessing different aspects of the RCM emulator), but **what metrics should we compute?**

### Distributional-based metrics

The purpose of this set of metrics is to assess **how well the emulator reproduces the distribution of the RCM**



**Source:** Baño-Medina, J., Manzanas, R., Cimadevilla, E., Fernández, J., González-Abad, J., Cofiño, A. S. & Gutiérrez, J. M. (2022). Downscaling multi-model climate projection ensembles with deep learning (DeepESD): contribution to CORDEX EUR-44. Geoscientific Model Development Discussions, 2022, 1-14.

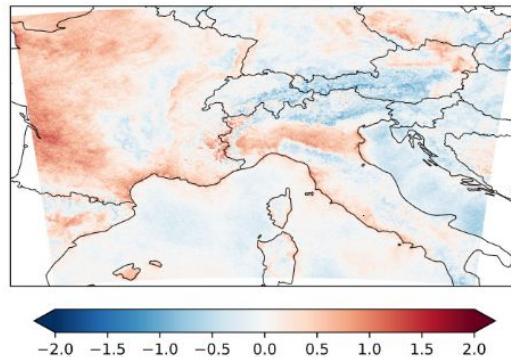
Ground-truth: RCM  
Prediction: RCM emulator

## Distributional-based Metrics

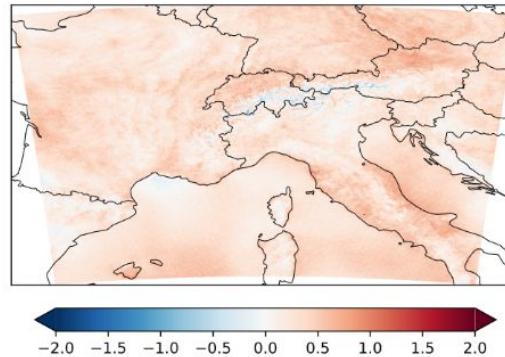
So far, we have reviewed the evaluation frameworks (i.e., what data to use for assessing different aspects of the RCM emulator), but **what metrics should we compute?**

### Distributional-based metrics

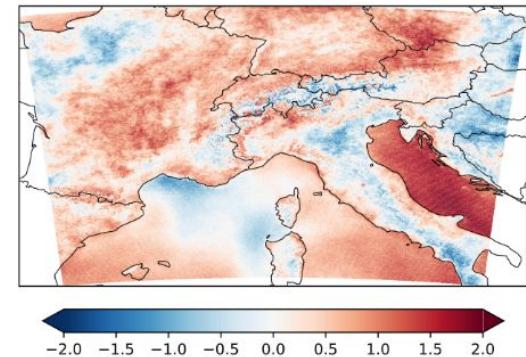
Bias in 2nd Percentile  
( $^{\circ}\text{C}$ )



Bias in mean  
( $^{\circ}\text{C}$ )



Bias in 98th Percentile  
( $^{\circ}\text{C}$ )



Variable: Temperature

Ground-truth: RCM

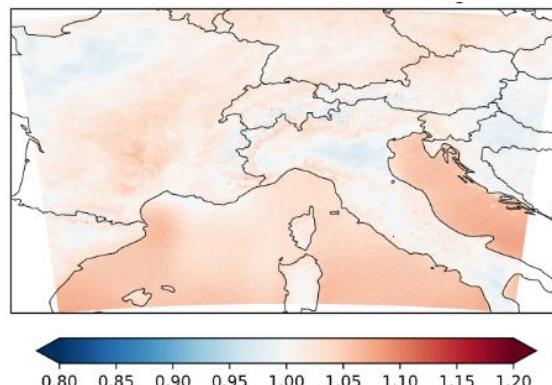
Prediction: RCM emulator

## Distributional-based Metrics

So far, we have reviewed the evaluation frameworks (i.e., what data to use for assessing different aspects of the RCM emulator), but **what metrics should we compute?**

### Distributional-based metrics

#### Ratio of Standard Deviations



Variable: Temperature  
Ground-truth: RCM  
Prediction: RCM emulator

## Example of Evaluation

Let us suppose that we have trained an RCM emulator on an RCM run over **GCM\_1** ([low emissions scenario](#))

## Example of Evaluation

Let us suppose that we have trained an RCM emulator on an RCM run over **GCM\_1** (**low emissions scenario**)

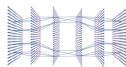
To evaluate the emulator's ability to simulate a **high emissions** scenario, I could proceed as follows  
**(soft-transferability)**

## Example of Evaluation

Let us suppose that we have trained an RCM emulator on an RCM run over **GCM\_1** (**low emissions scenario**)

To evaluate the emulator's ability to simulate a **high emissions** scenario, I could proceed as follows  
**(soft-transferability)**

RCM  
emulator

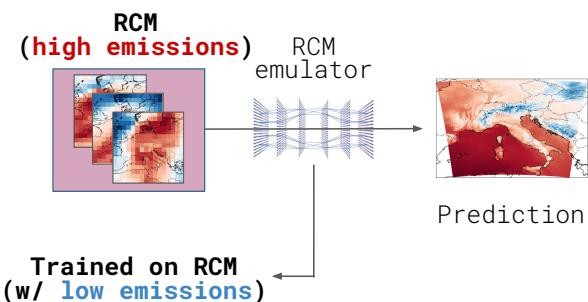


Trained on RCM  
(w/ **low emissions**)

## Example of Evaluation

Let us suppose that we have trained an RCM emulator on an RCM run over **GCM\_1** (**low emissions scenario**)

To evaluate the emulator's ability to simulate a **high emissions** scenario, I could proceed as follows  
**(soft-transferability)**



## Example of Evaluation

Let us suppose that we have trained an RCM emulator on an RCM run over **GCM\_1** (**low emissions scenario**)

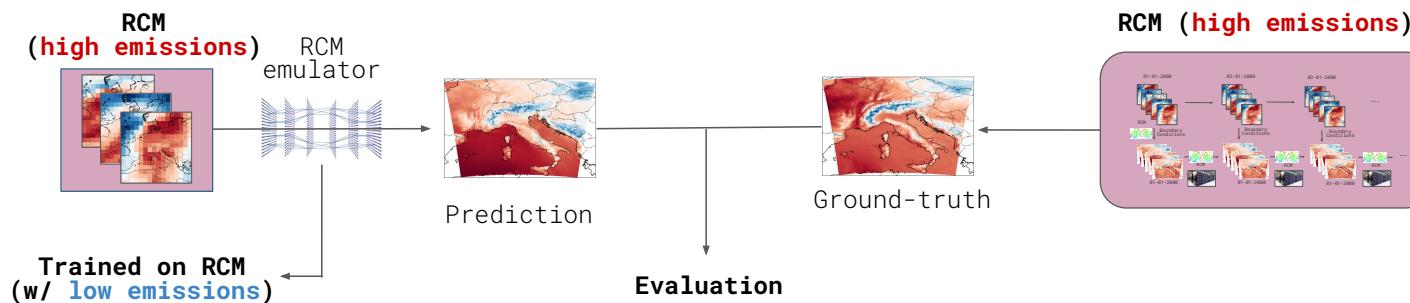
To evaluate the emulator's ability to simulate a **high emissions** scenario, I could proceed as follows  
**(soft-transferability)**



## Example of Evaluation

Let us suppose that we have trained an RCM emulator on an RCM run over **GCM\_1** (**low emissions scenario**)

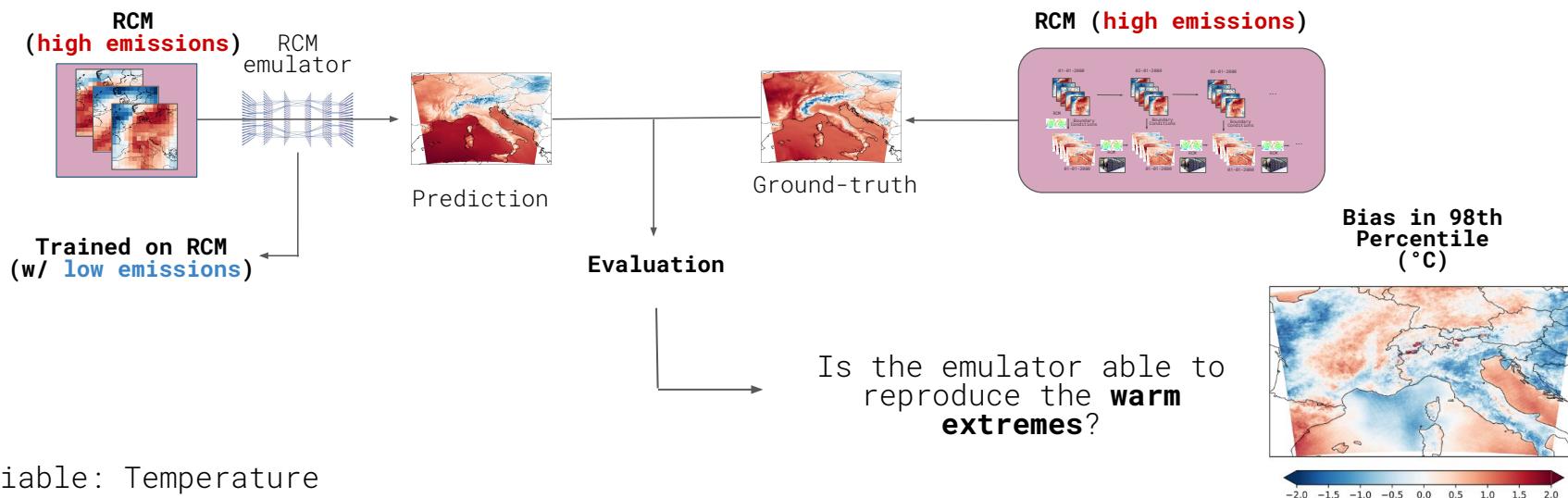
To evaluate the emulator's ability to simulate a **high emissions** scenario, I could proceed as follows  
**(soft-transferability)**



## Example of Evaluation

Let us suppose that we have trained an RCM emulator on an RCM run over **GCM\_1** (**low emissions scenario**)

To evaluate the emulator's ability to simulate a **high emissions** scenario, I could proceed as follows  
**(soft-transferability)**



# Table of Contents

## Part I: Review

- DL for Weather/Climate
- RCM Emulation
- Conclusions

## Part II: Training a RCM emulator

- Problem Statement
- Training Frameworks
- Training and Inference

## Part III: Evaluating a RCM emulator

- Importance of Proper Evaluation
- Soft-Transferability
- Hard-Transferability
- Metrics

## Part IV: The Importance of Benchmarks

- What is a benchmark?
- CORDEX-ML-Bench

## Benchmark

A benchmark in DL is a **standardized dataset** and **evaluation metric** used to objectively compare and **track the performance of different models**

# Benchmark

A benchmark in DL is a **standardized dataset** and **evaluation metric** used to objectively compare and **track the performance of different models**



**ImageNet** is a large-scale labeled image dataset widely used to train and benchmark deep learning models for **image classification**

# Benchmark

A benchmark in DL is a **standardized dataset** and **evaluation metric** used to objectively compare and **track the performance of different models**



**ImageNet** is a large-scale labeled image dataset widely used to train and benchmark deep learning models for **image classification**

There are multiple benchmarks for **LLMs** that try to **evaluate different dimensions**, such as coding, reasoning, or mathematical skills

## LiveBench

### A Challenging, Contamination-Free LLM Benchmark

LiveBench will appear as a [Spotlight Paper](#) in ICLR 2025.

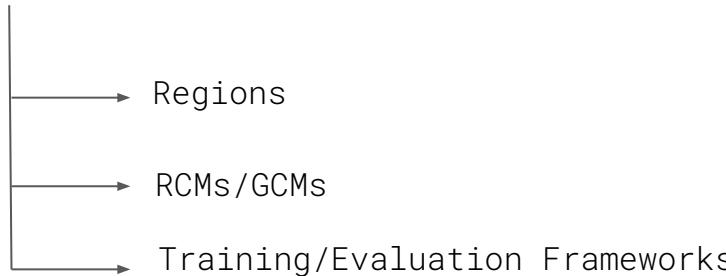
This work is sponsored by [Abacus.AI](#)



Model	Organization	Global Average	Reasoning Average	Coding Average	Agentic Coding Average
GPT-5 High	OpenAI	78.59	98.17	75.31	43.33
GPT-5 Medium	OpenAI	76.45	96.58	73.25	35.00
GPT-5 Low	OpenAI	75.34	90.47	72.49	41.67
o3 Pro High	OpenAI	74.72	94.67	76.78	31.67
o3 High	OpenAI	74.61	94.67	76.71	36.67
Claude 4.1 Opus Thinking	Anthropic	73.48	93.19	73.96	33.33

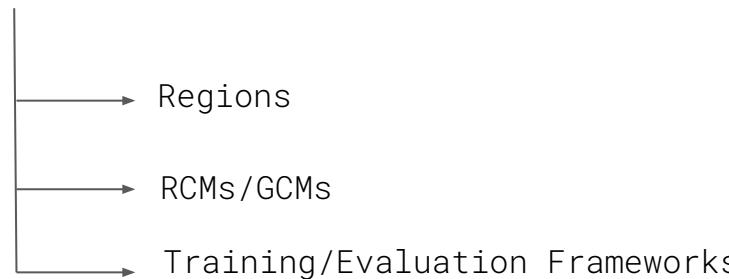
## Benchmarks in RCM emulation

Despite being a recent field, **various papers** have already been published on **RCM emulation** spanning **different** . . .

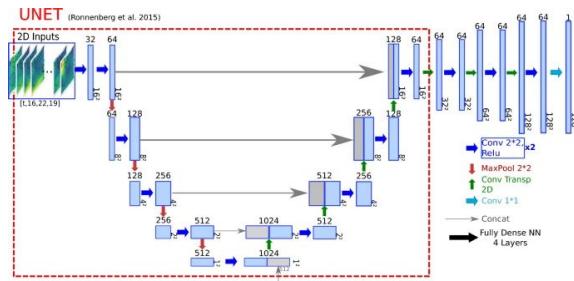


## Benchmarks in RCM emulation

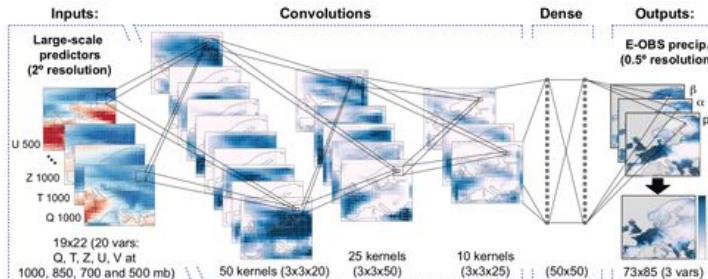
Despite being a recent field, **various papers** have already been published on **RCM emulation** spanning **different**...



This **complicates the intercomparison of different DL architectures** for RCM emulation



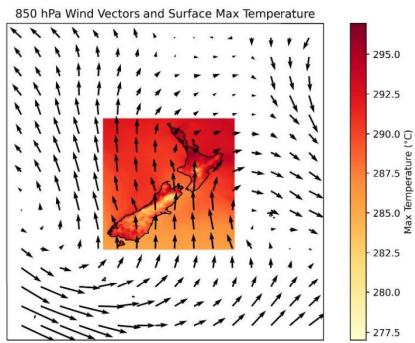
Source: Doury, A., Somot, S., & Gadat, S. (2024). On the suitability of a convolutional neural network based RCM-emulator for fine spatio-temporal precipitation. *Climate Dynamics*, 62(9), 8587-8613.



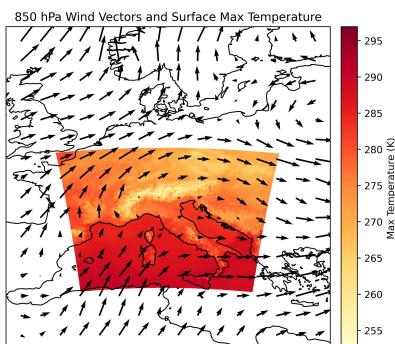
Source: Baño-Medina, J., Manzanas, R., & Gutiérrez, J. M. (2020). Configuration and intercomparison of deep learning neural models for statistical downscaling. *Geoscientific Model Development*, 13(4), 2109-2124.

## What is CORDEX-ML-Bench?

Within CORDEX we are developing **CORDEX-ML-Bench**, a benchmarking dataset for **AI-based regional climate downscaling**



Different regions



Different domains

Different training experiments

Different evaluation experiments



## What is CORDEX-ML-Bench?

In the next session, we will introduce a **first version of this benchmark** (still a **work in progress**) and demonstrate **how to train and evaluate RCM emulators**

## What is CORDEX-ML-Bench?

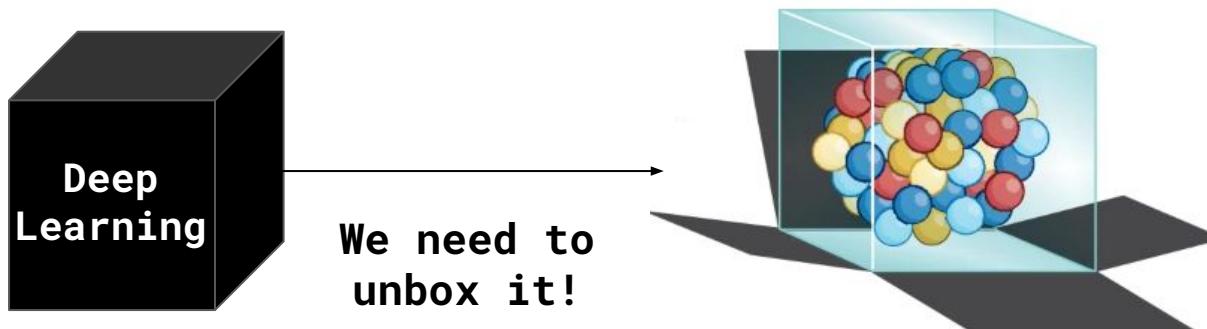
In the next session, we will introduce a **first version of this benchmark** (still a **work in progress**) and demonstrate **how to train and evaluate RCM emulators**

## Questions?

# Interpretability

Despite being formed by simple elements, its **composition** makes deep learning models to be considered **black boxes**

This makes it difficult to gain a **comprehensive** understanding of their inner functioning, **particularly for downscaling climate change projections**



# Interpretability

**Interpretability techniques** emerged in the computer vision field to explain the functioning and results of deep learning models

Brushing teeth



Cutting trees

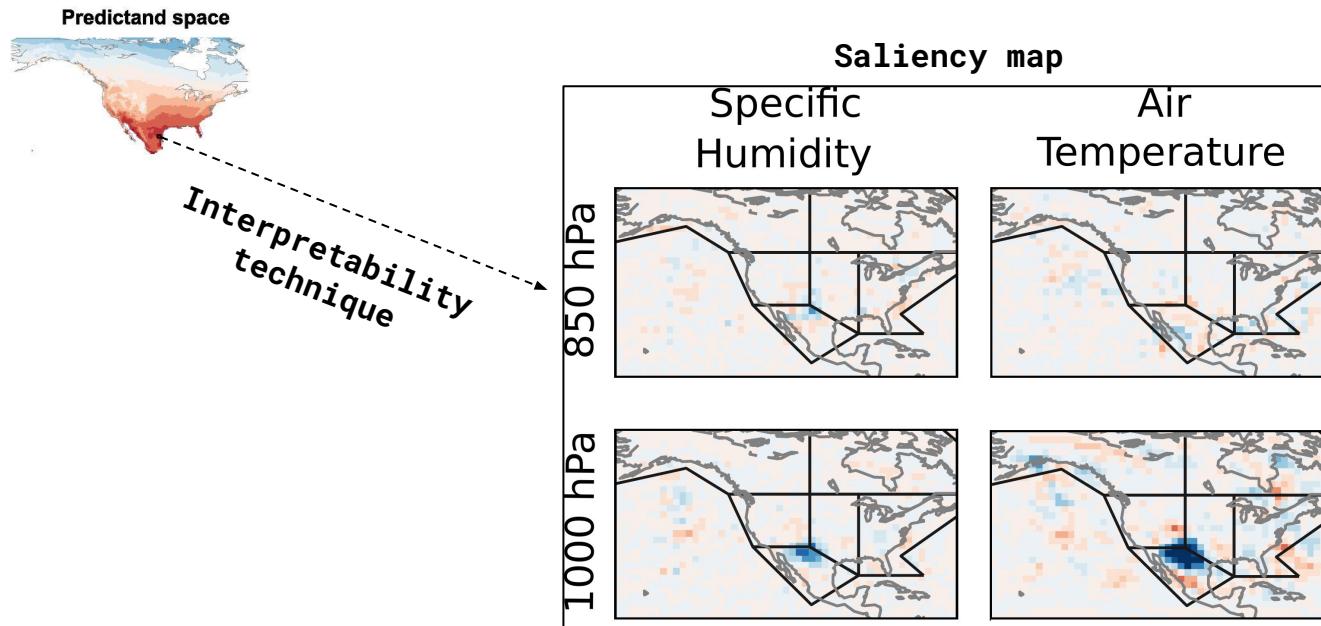


**Saliency maps** assign to each feature a **relevance score** representing its influence on the computed prediction

Source: Zhou, Bolei, et al. "Learning deep features for discriminative localization." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016

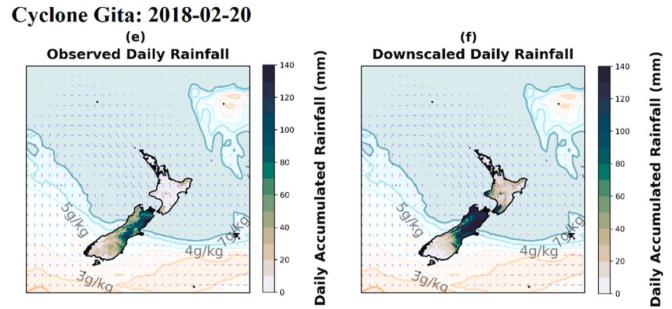
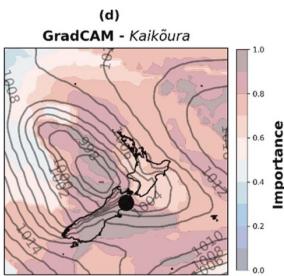
# Interpretability

For **statistical downscaling**, **saliency maps** are defined over the **grid-points of the large-scale variables** conforming the **predictor**



# Interpretability

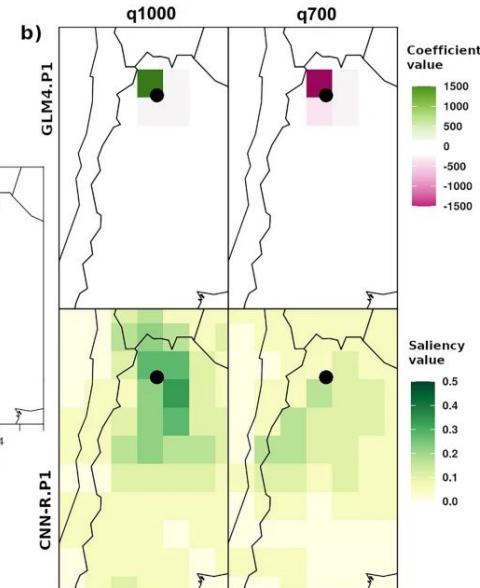
**Interpretability techniques** have been used to **gain insights into the relationships learnt by the deep learning model**



Rampal et al. (2022)

The deep learning model is taking into account the **cyclone** to downscale the **precipitation extreme**

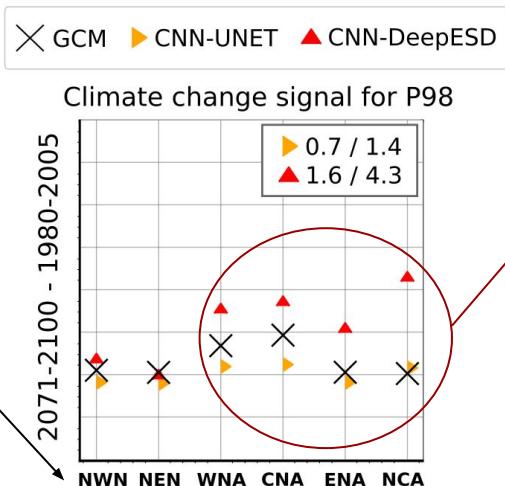
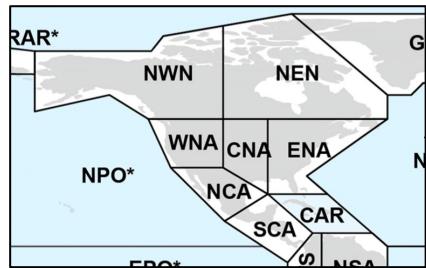
Balmaceda-Huarte et al. (2024)



The deep learning model is **learning suitable relationships**

# Interpretability

**Interpretability techniques** have even been used to **debug** deep learning models



**DeepESD** is learning a **spurious pattern for southern regions**, hampering its ability to extrapolate

When downscaling the GCM (**extrapolation**) the DeepESD overestimates the climate change signal for the southern regions of North America

