



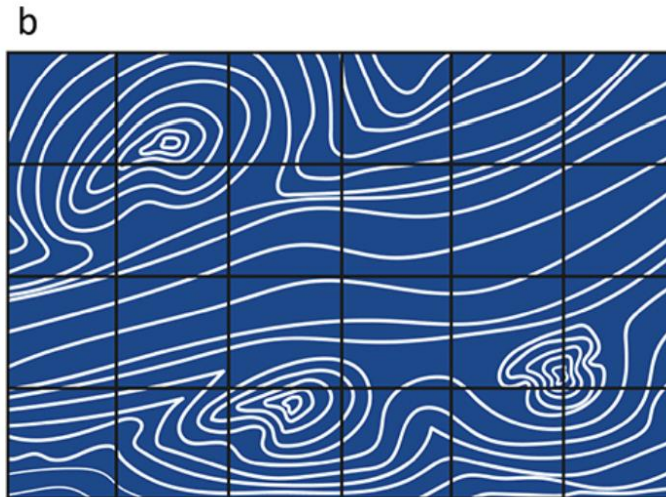
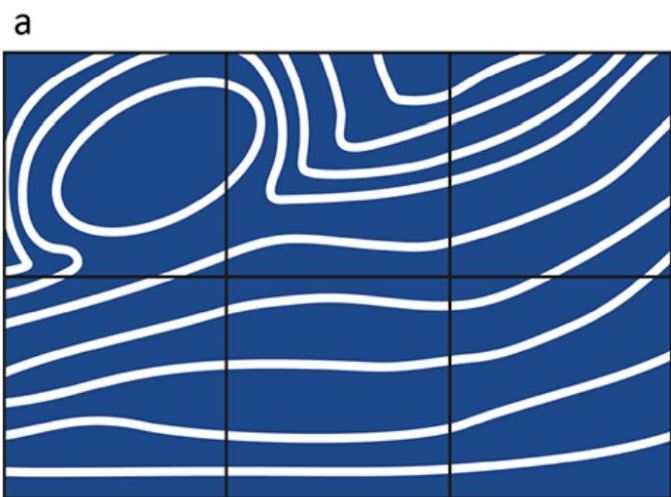
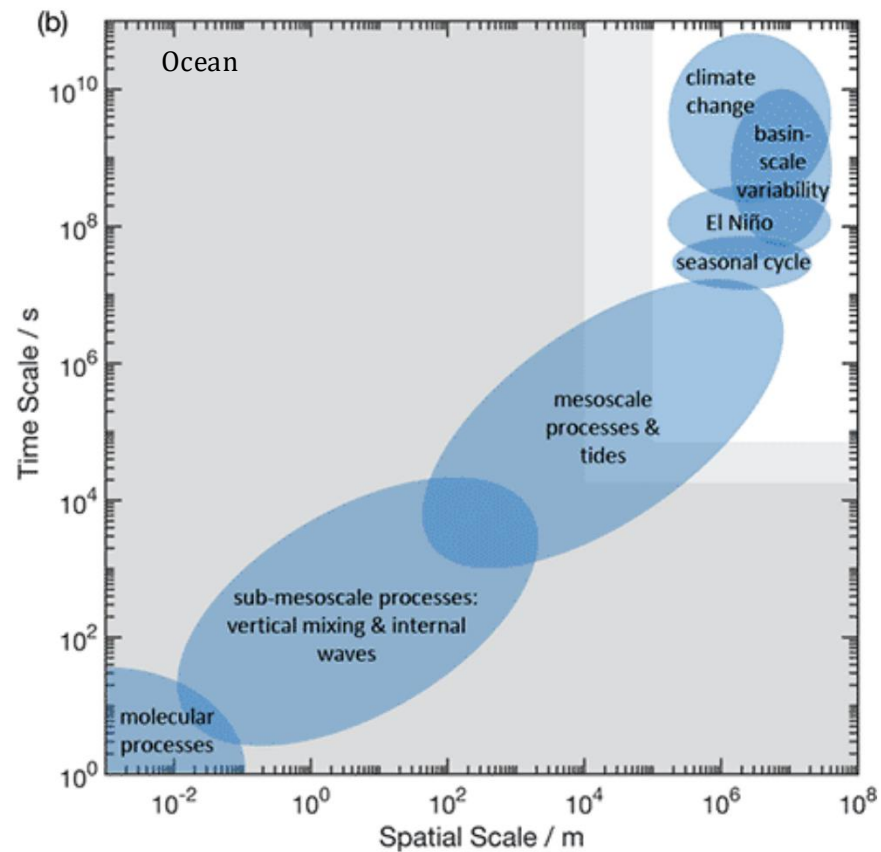
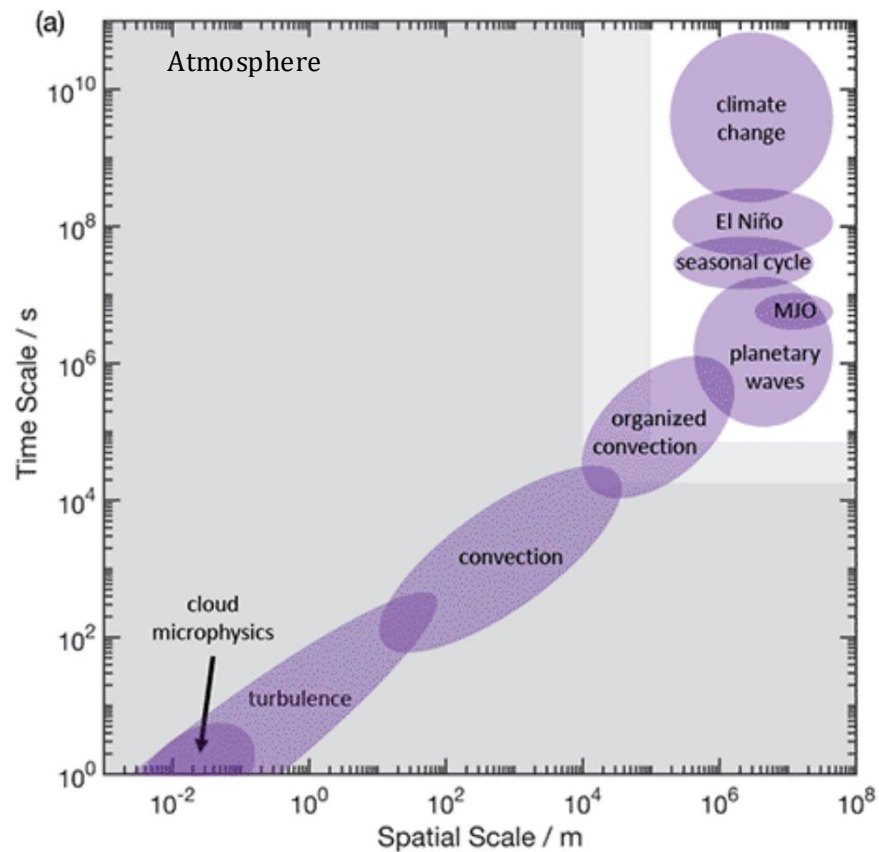
# **Synthesizing data sources**

# **Data assimilation**

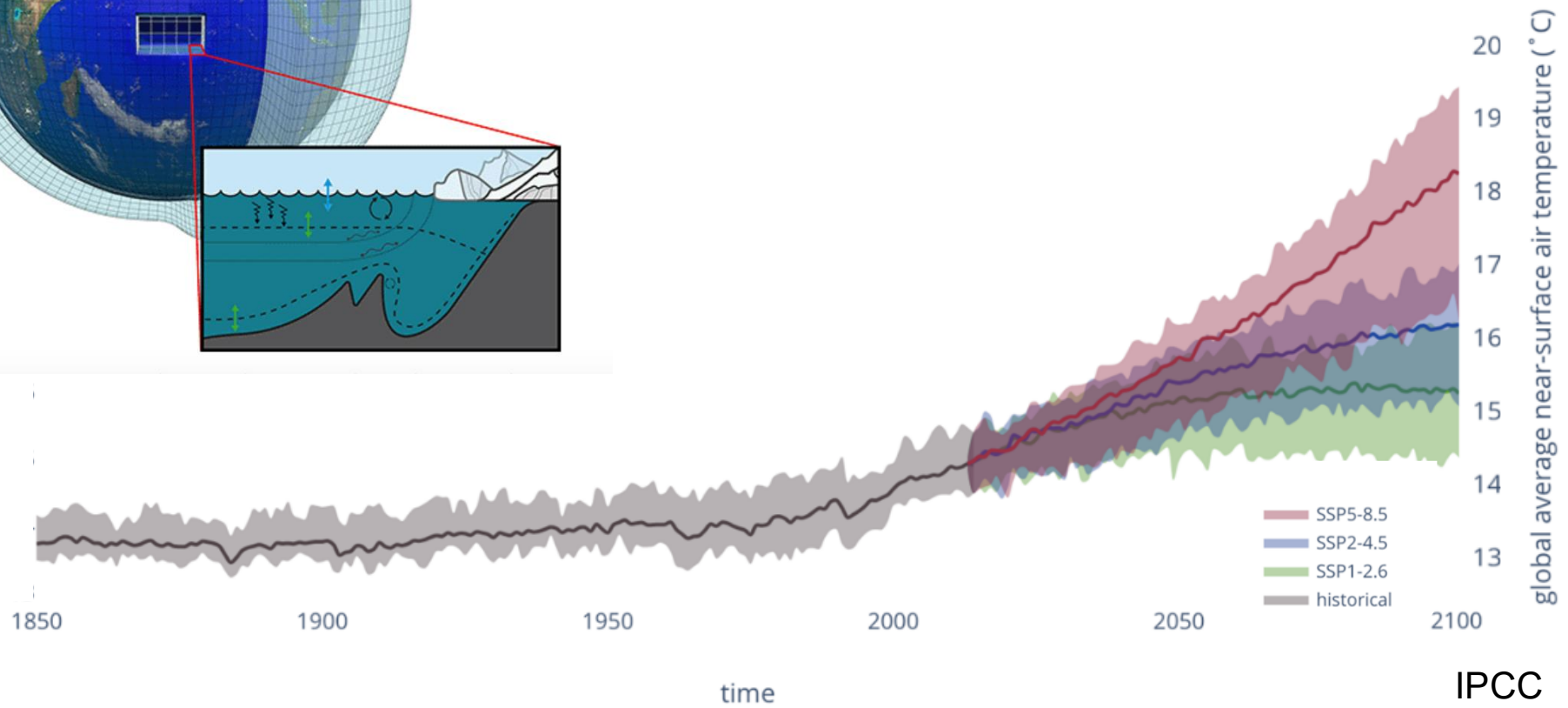
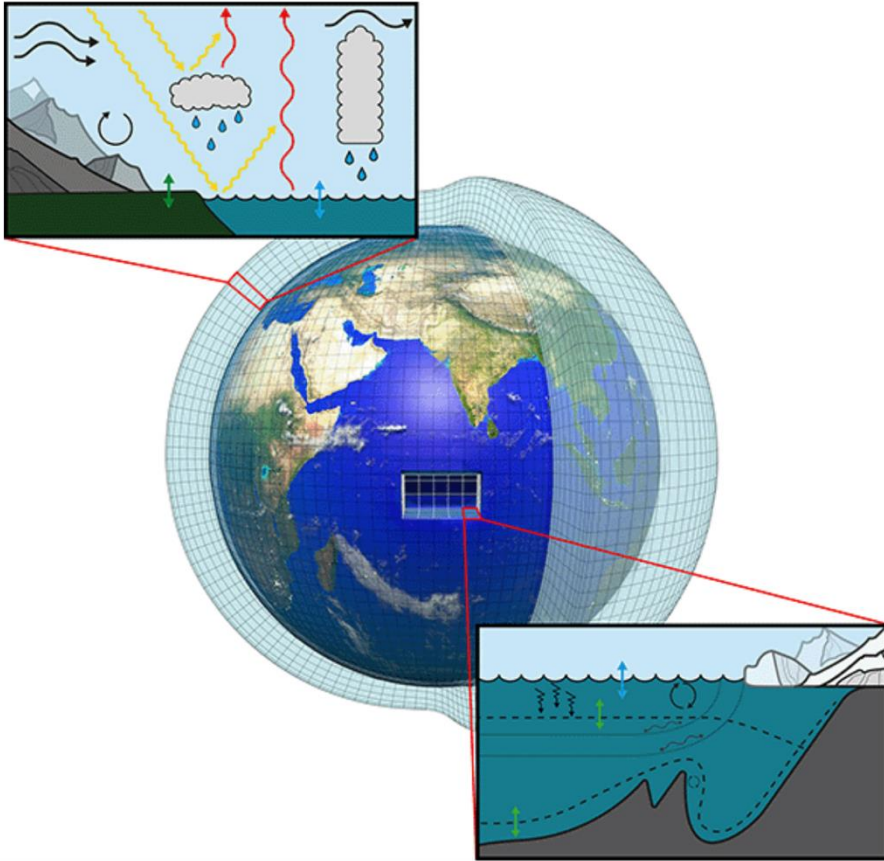
**12.S992 AI for Climate Action**

**Spring 2026**

**Speaker: Abigail Bodner**



# Global Climate Change

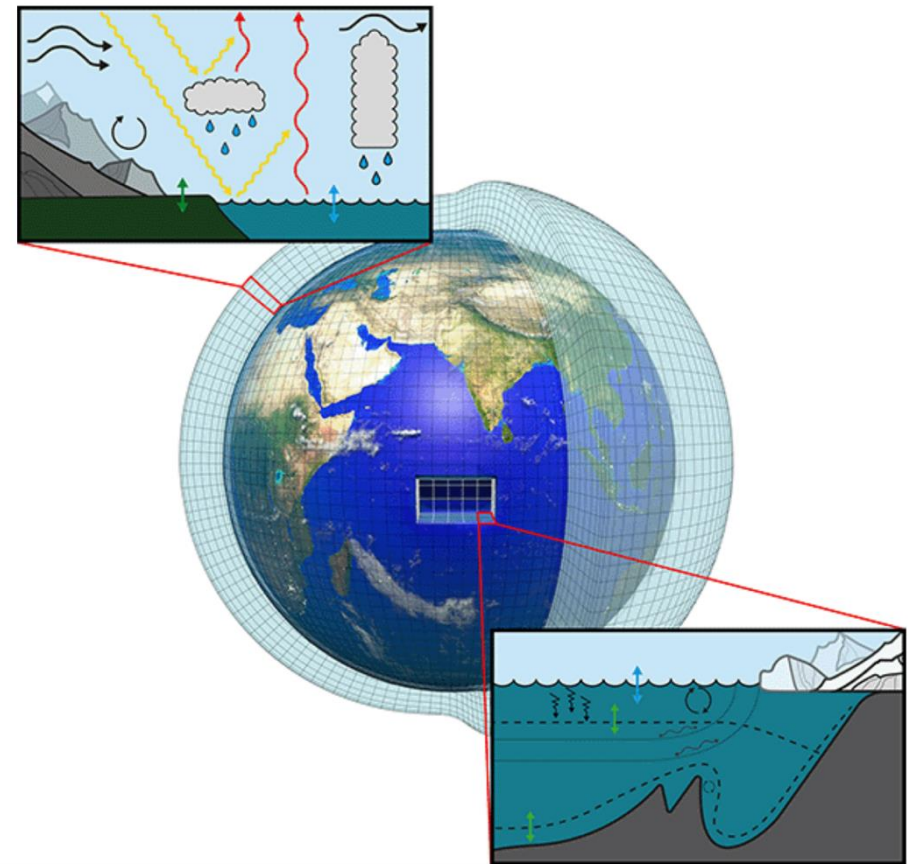


IPCC

# Data assimilation (DA)

Also called reanalysis or state estimate

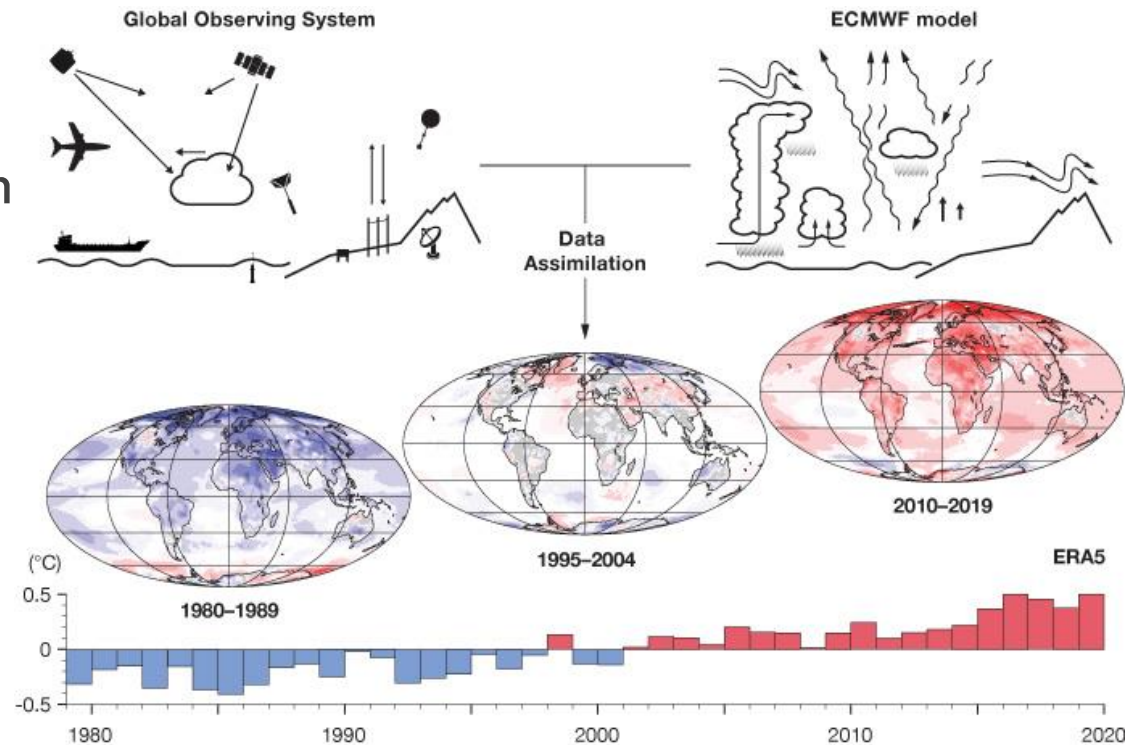
Unlike the parameterization problem, which solves the equations of motion and infers subgrid effects, in DA the goal is to create a gridded product from all available data sources.



# Data assimilation (DA)

Also called reanalysis or state estimate

Unlike the parameterization problem, which solves the equations of motion and infers subgrid effects, in DA the goal is to create a gridded product from all available data sources.

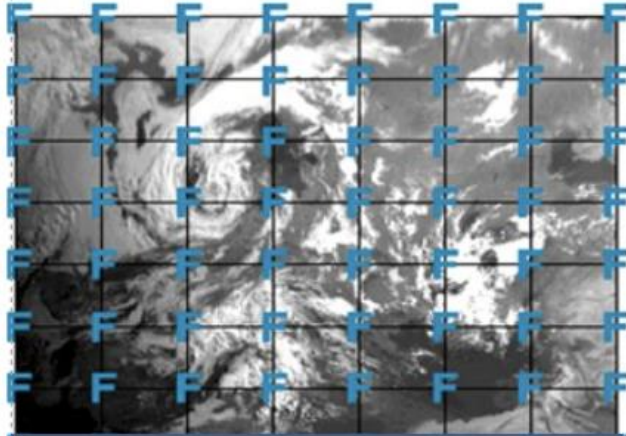


# Data assimilation (DA)

Quantify the state of the system (ocean/atmosphere/land/cryosphere) given all available information from:

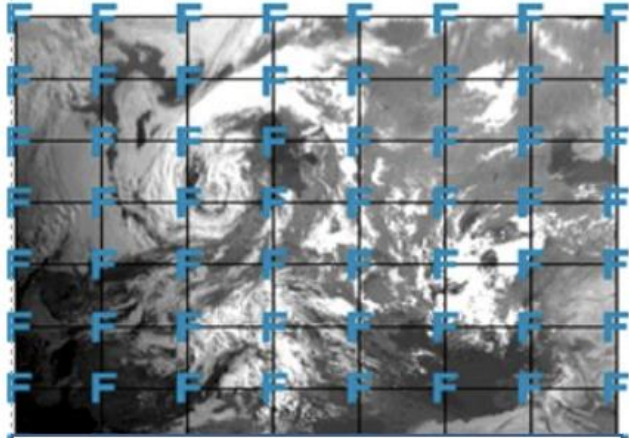
- Models
- Observations

## The data assimilation process

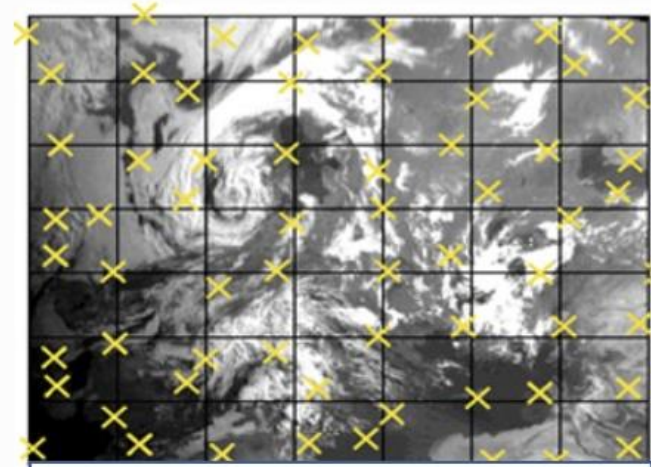


Model forecast (with errors)

## The data assimilation process

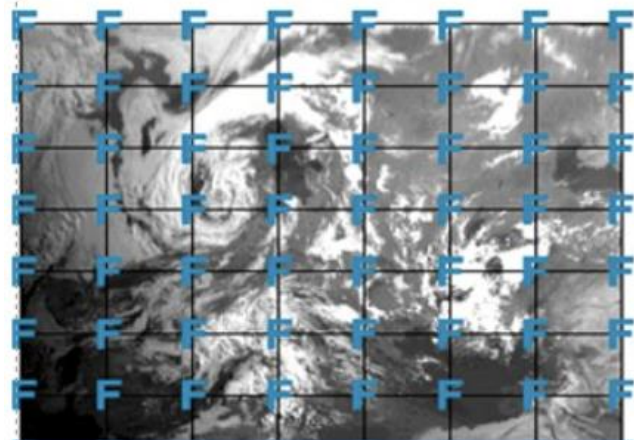


Model forecast (with errors)

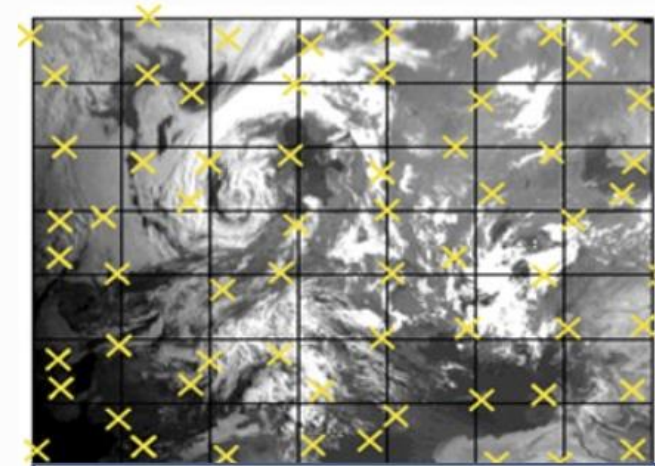


Observations (with errors)

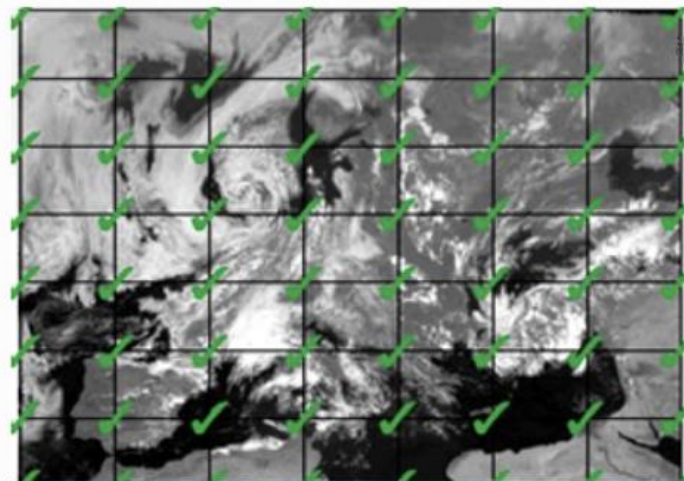
## The data assimilation process



Model forecast (with errors)



Observations (with errors)



Analysis (with smaller errors)

# Data assimilation (DA)

Quantify the state of the system (ocean/atmosphere/land/cryosphere) given all available information from:

- Models
- Observations

Challenges:

# Data assimilation (DA)

Quantify the state of the system (ocean/atmosphere/land/cryosphere) given all available information from:

- Models
- Observations

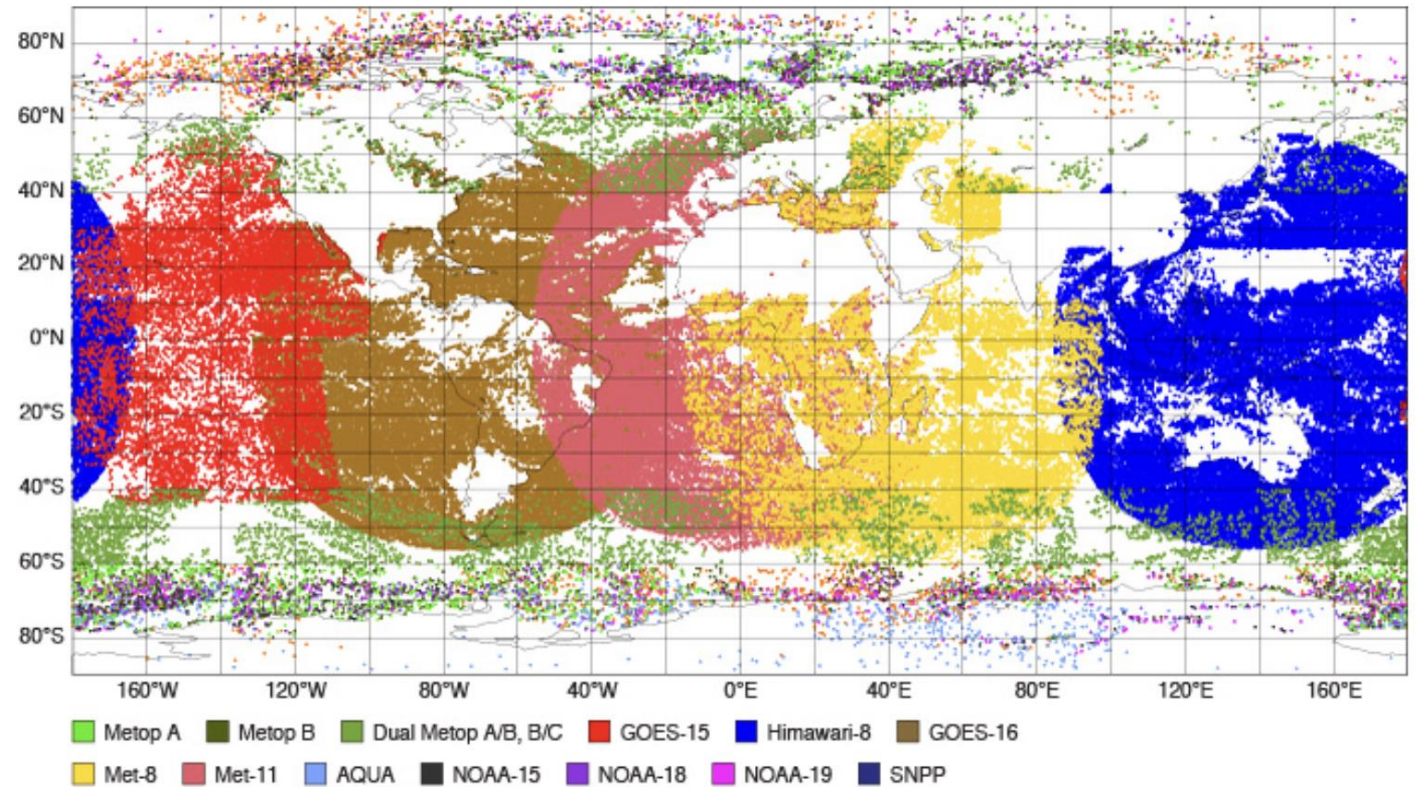
Challenges:

- What is the best data (error estimates)
- How to synthesize data from different sources, i.e. time, space, resolution, etc.

# Example: Forecasting (4D-Var)

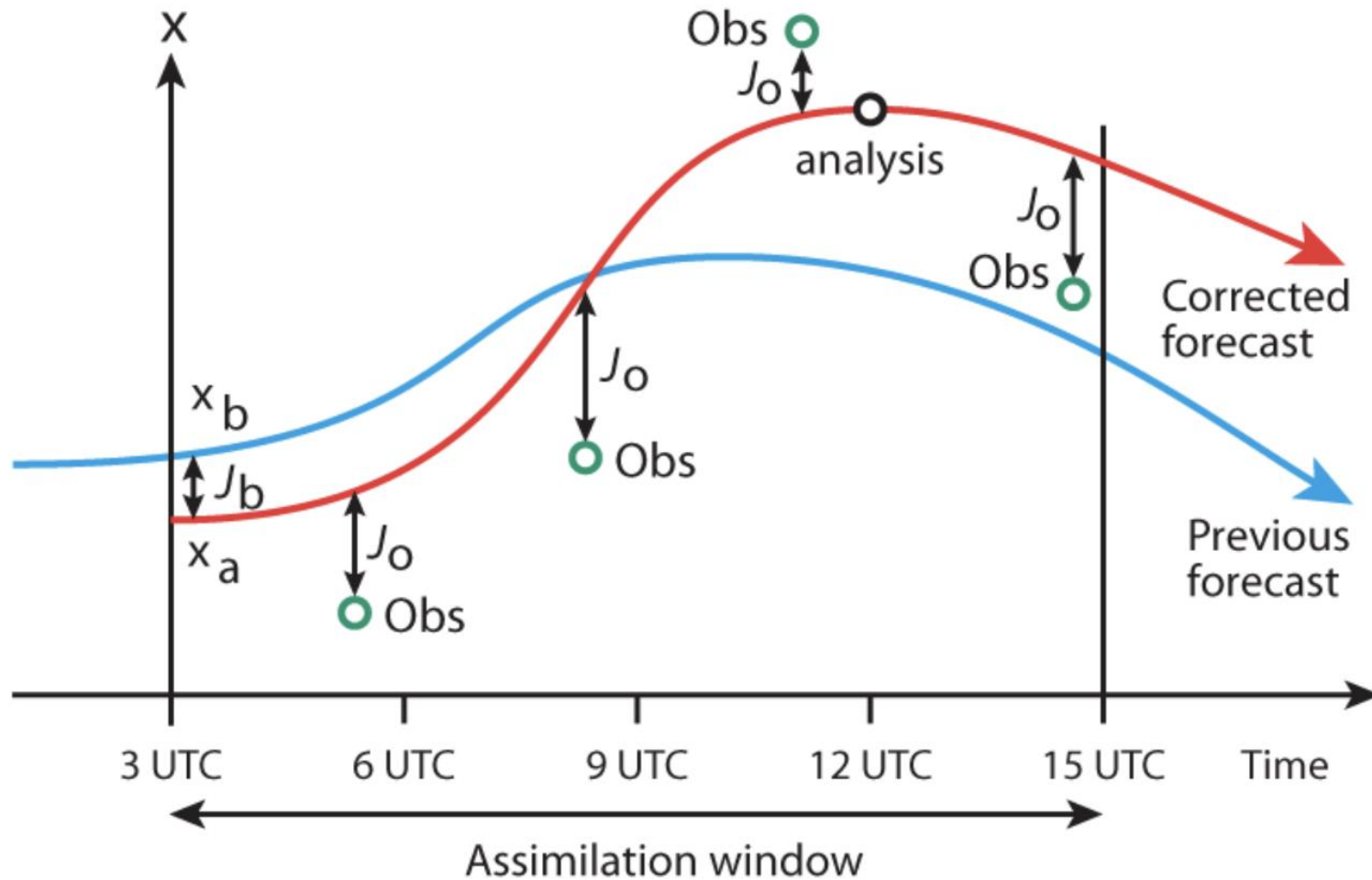


Weather observations come from many sources, but they cannot provide a complete picture of the state of the Earth system at a given point in time.  
(Diagram: WMO)



Typical coverage of active atmospheric motion vector (AMV) data for a 12-hour assimilation cycle (00 UTC, 7 March 2019).

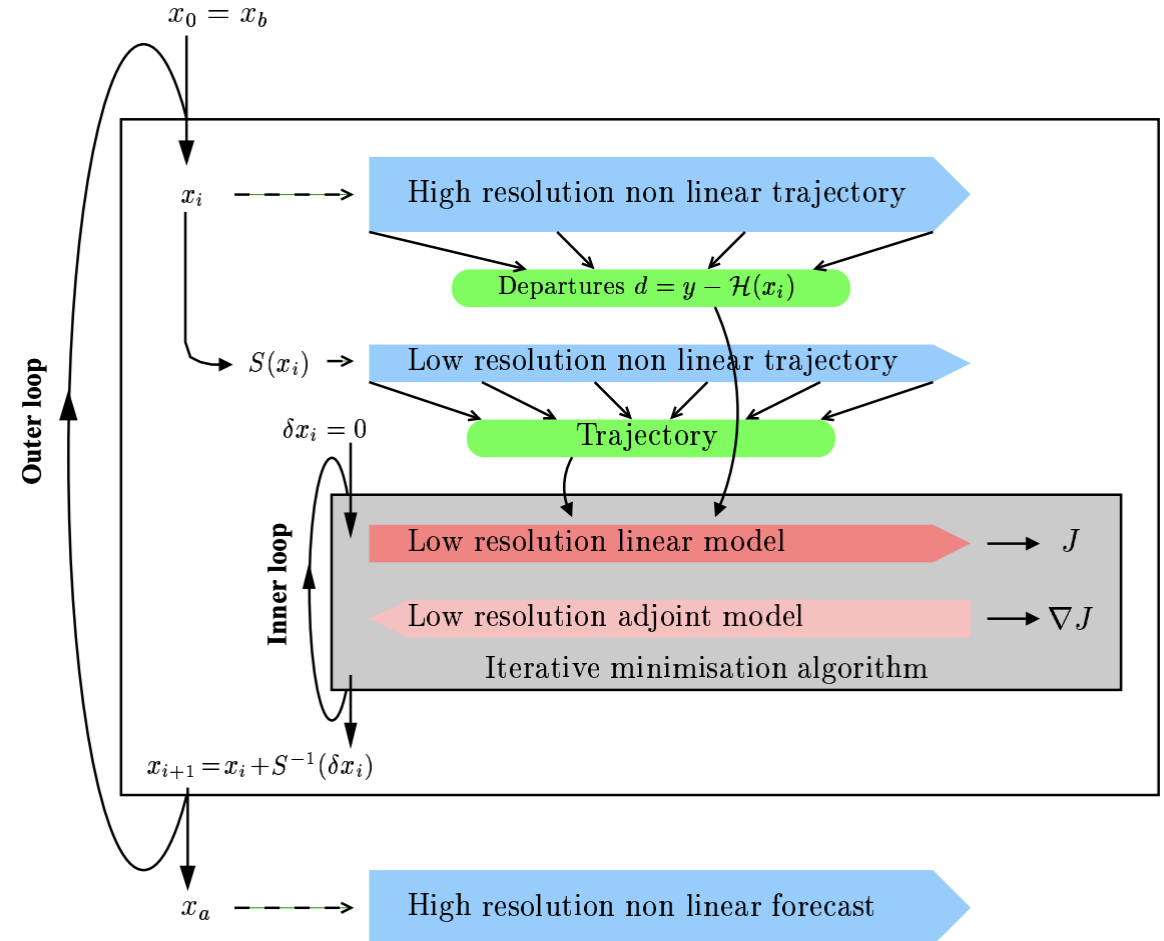
# Example: Forecasting (4D-Var)



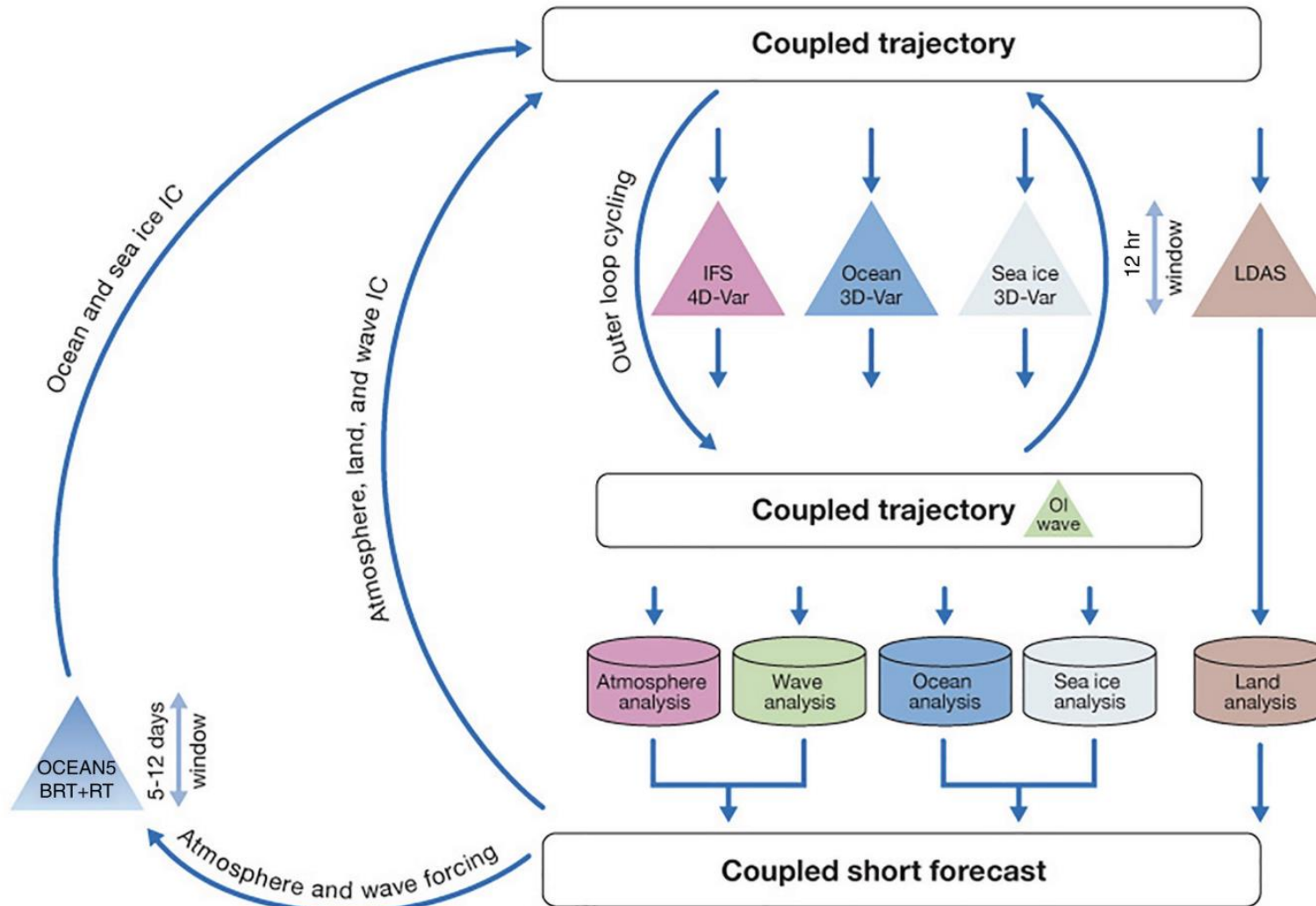
# Example: Forecasting (4D-Var)

$$J(\mathbf{x}) = \frac{1}{2} (\mathbf{x} - \mathbf{x}_b)^T \mathbf{B}^{-1} (\mathbf{x} - \mathbf{x}_b) + \frac{1}{2} \sum_k (\mathcal{H}_k \mathcal{M}_k(\mathbf{x}) - \mathbf{y}_k)^T \mathbf{R}_k^{-1} (\mathcal{H}_k \mathcal{M}_k(\mathbf{x}) - \mathbf{y}_k)$$

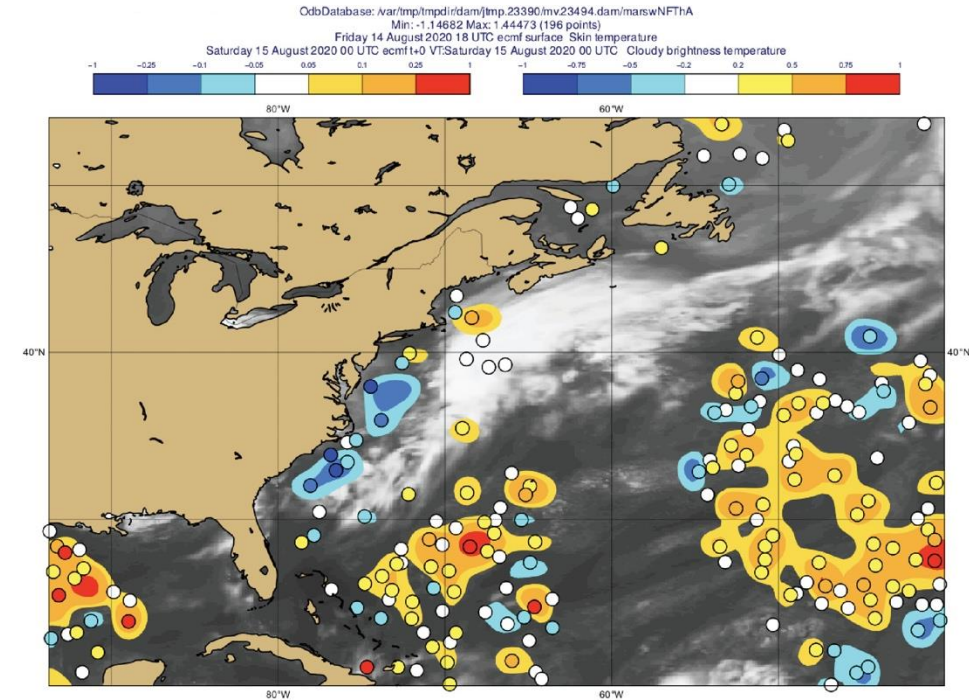
$$\nabla J(\mathbf{x}) = \mathbf{B}^{-1} (\mathbf{x} - \mathbf{x}_b) + \sum_k \mathbf{M}_k^T \mathbf{H}_k^T \mathbf{R}_k^{-1} (\mathcal{H}_k \mathcal{M}_k(\mathbf{x}) - \mathbf{y}_k).$$



# Example: Forecasting (4D-Var)

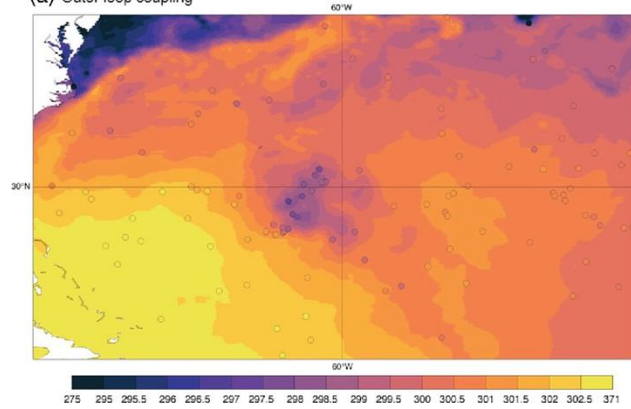


# Example: Forecasting (4D-Var)

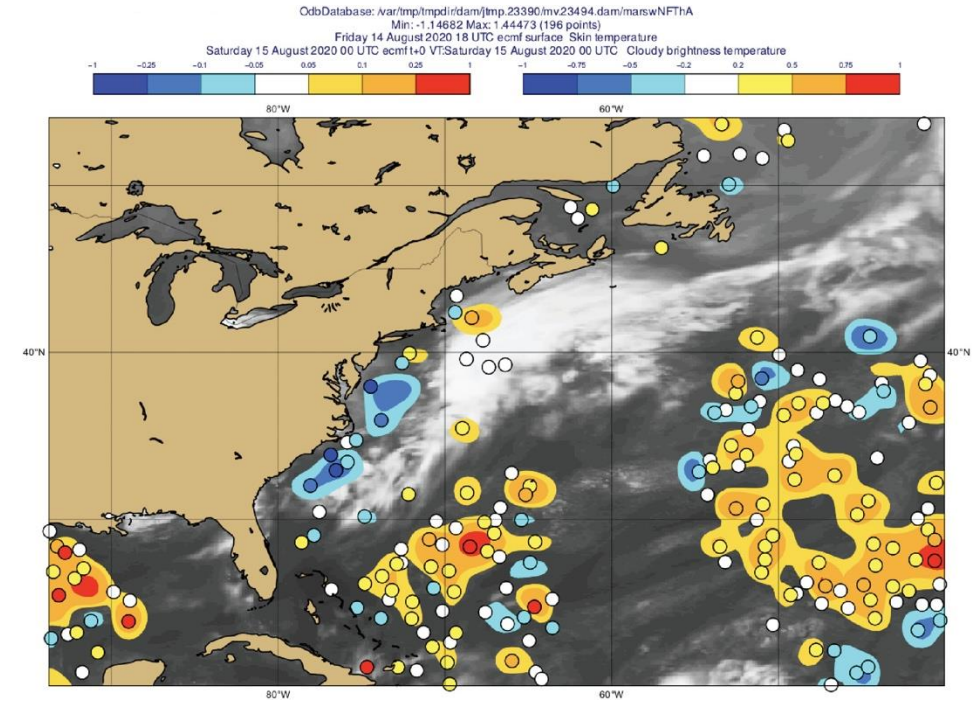
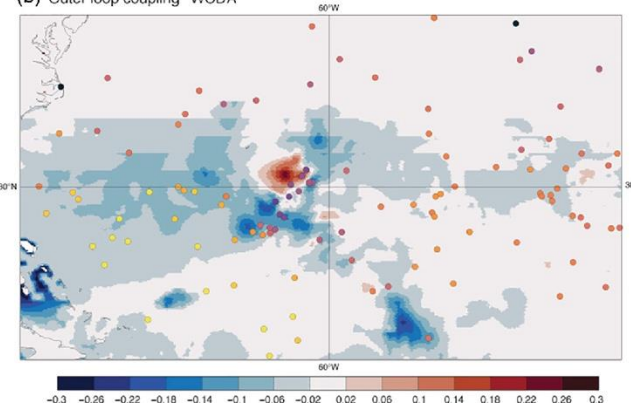


# Example: Forecasting (4D-Var)

(a) Outer loop coupling

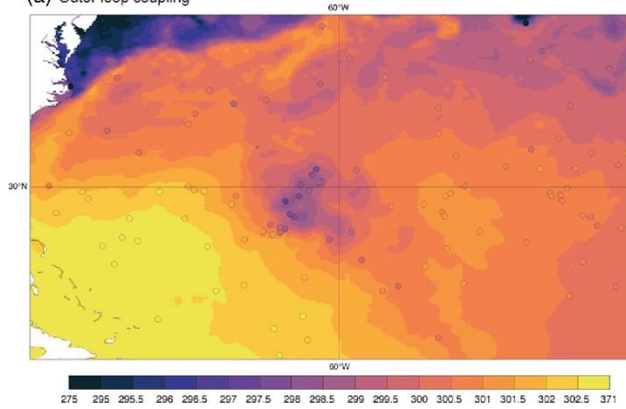


(b) Outer loop coupling-WCDA

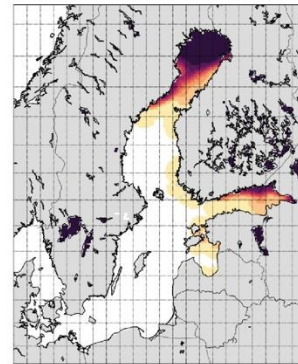
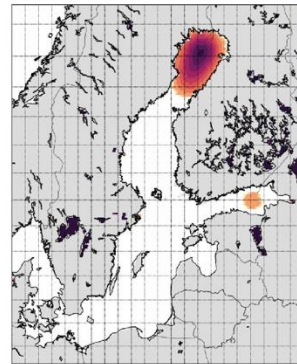
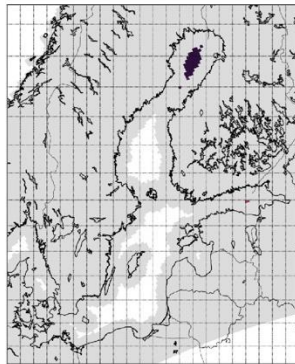
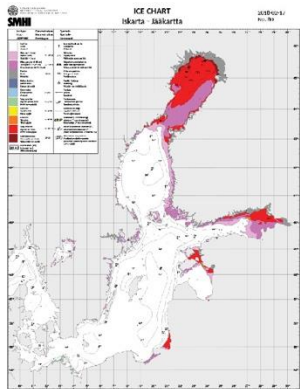
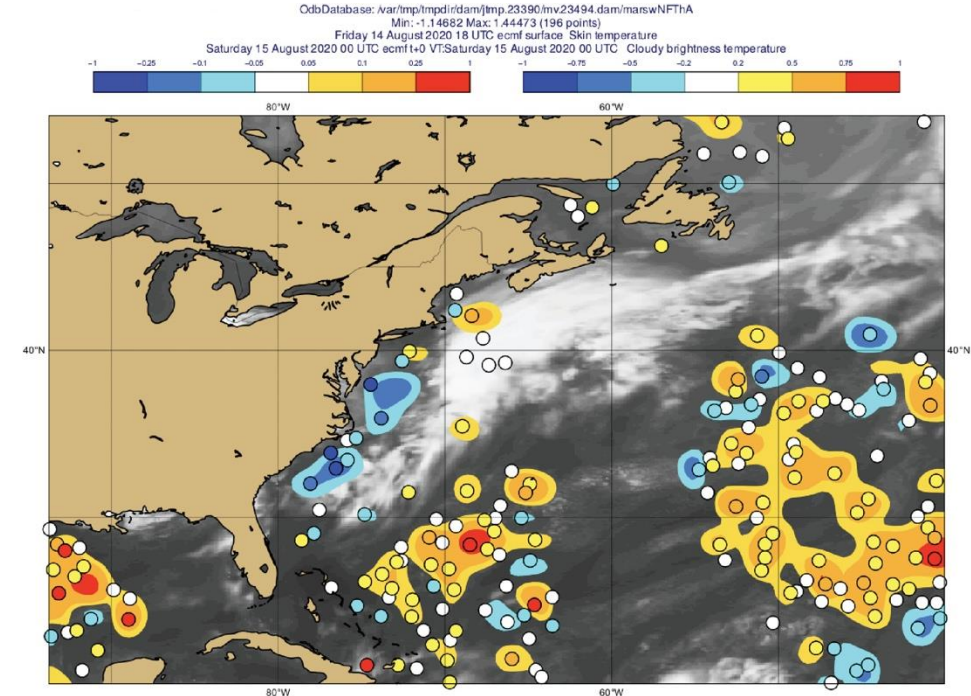
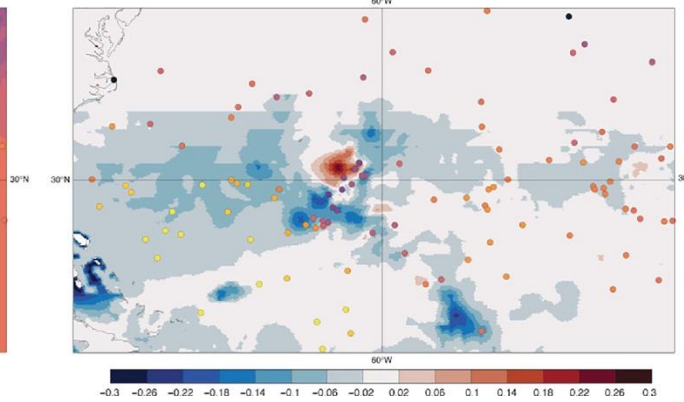


# Example: Forecasting (4D-Var)

(a) Outer loop coupling



(b) Outer loop coupling-WCDA

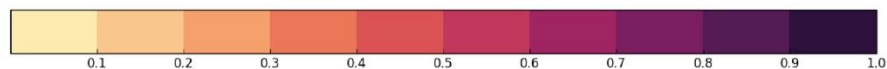


(a) Manual ice chart

(b) L3 from OSI SAF

(c) Uncoupled analysis

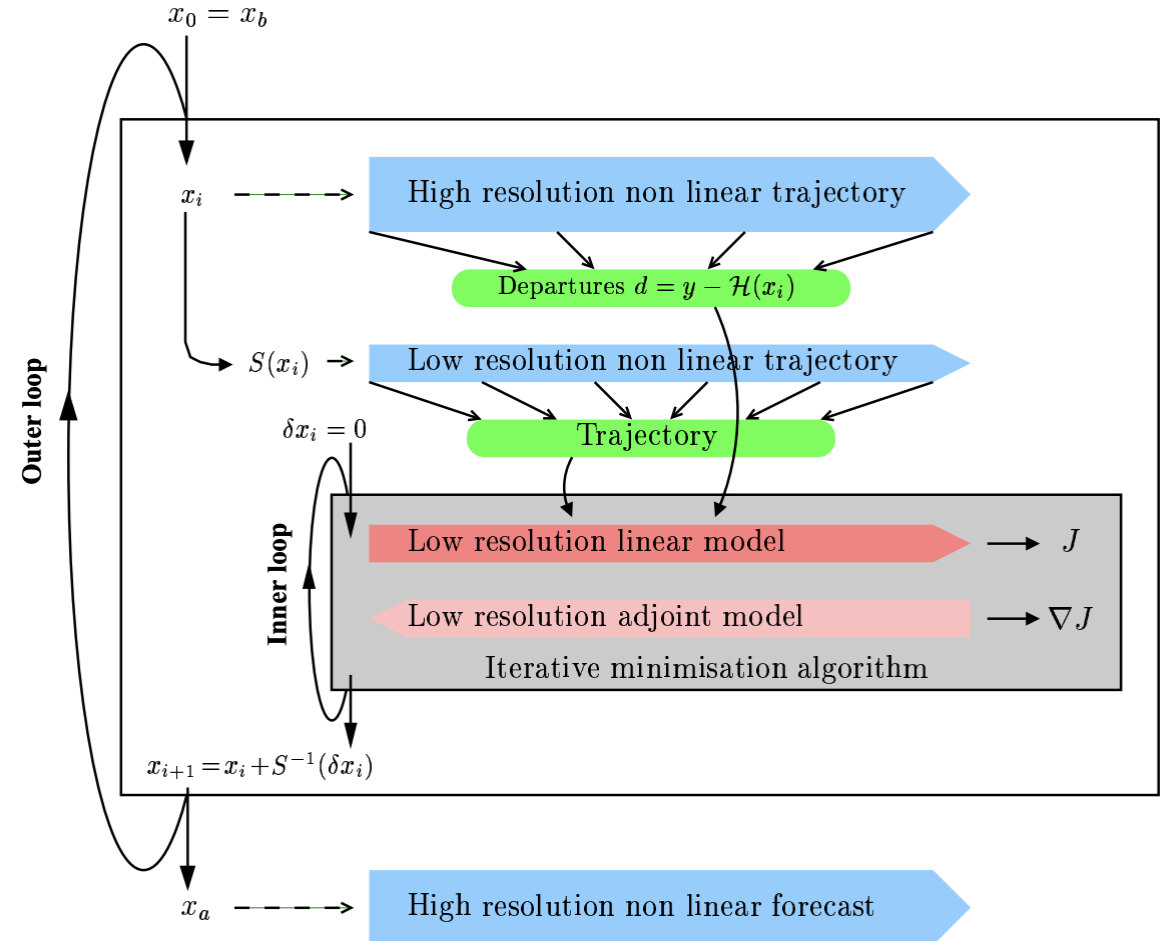
(d) WCDA



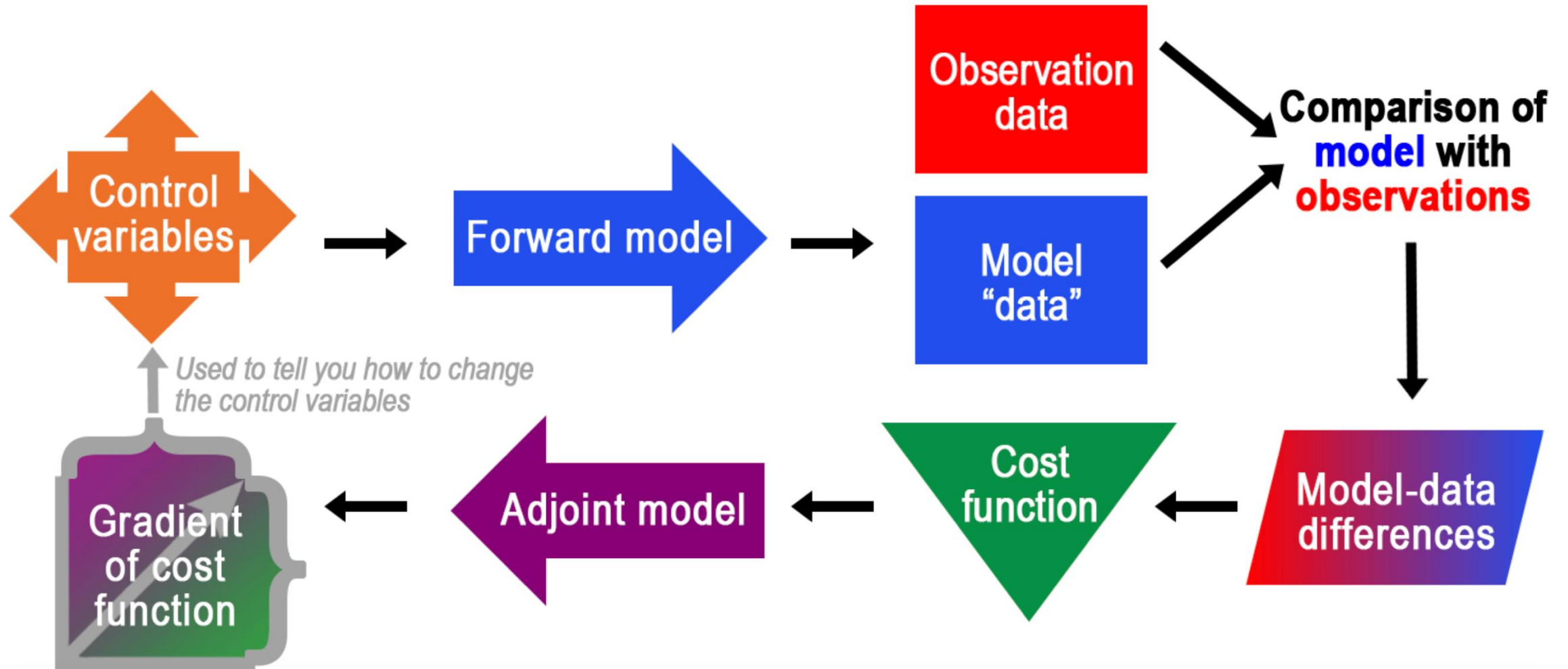
# Example: Forecasting (4D-Var)

$$J(\mathbf{x}) = \frac{1}{2} (\mathbf{x} - \mathbf{x}_b)^T \mathbf{B}^{-1} (\mathbf{x} - \mathbf{x}_b) + \frac{1}{2} \sum_k (\mathcal{H}_k \mathcal{M}_k(\mathbf{x}) - \mathbf{y}_k)^T \mathbf{R}_k^{-1} (\mathcal{H}_k \mathcal{M}_k(\mathbf{x}) - \mathbf{y}_k)$$

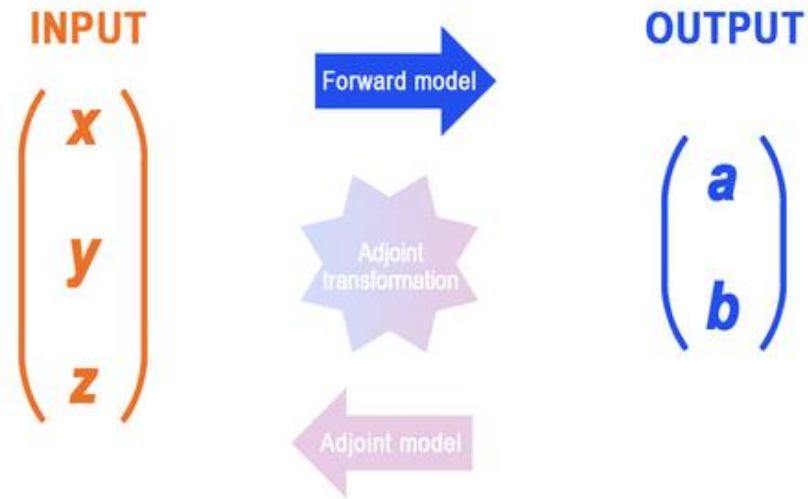
$$\nabla J(\mathbf{x}) = \mathbf{B}^{-1} (\mathbf{x} - \mathbf{x}_b) + \sum_k \mathbf{M}_k^T \mathbf{H}_k^T \mathbf{R}_k^{-1} (\mathcal{H}_k \mathcal{M}_k(\mathbf{x}) - \mathbf{y}_k).$$



# Example: State estimate (3D-Var)



# Example: State estimate (3D-Var)



$$a = x - 2y + 3z$$

$$\begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 1 & -2 & 3 \\ 4 & 0 & -5 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix}$$

Here's the other equation...  $b = 4x - 5z$

Forward model:

$$\begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 1 & -2 & 3 \\ 4 & 0 & -5 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix}$$

Adjoint model:

$$\begin{pmatrix} \hat{x} \\ \hat{y} \\ \hat{z} \end{pmatrix} = \begin{pmatrix} 1 & 4 \\ -2 & 0 \\ 3 & -5 \end{pmatrix} \begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix}$$

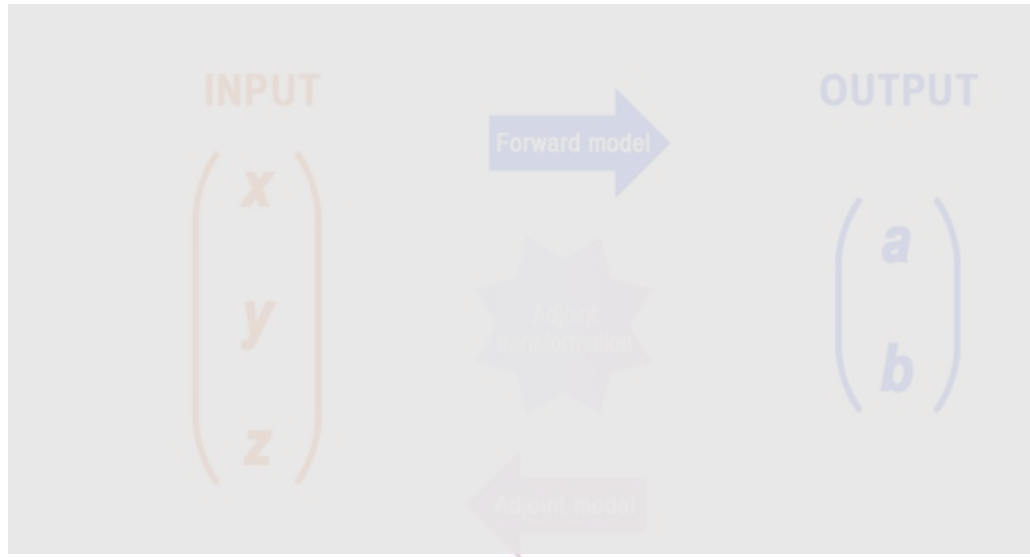
$$\begin{pmatrix} \hat{x} \\ \hat{y} \\ \hat{z} \end{pmatrix} = \begin{pmatrix} 1 & 4 \\ -2 & 0 \\ 3 & -5 \end{pmatrix} \begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix}$$

$$\hat{x} = \hat{a} + 4\hat{b}$$

$$\hat{y} = -2\hat{a}$$

$$\hat{z} = 3\hat{a} - 5\hat{b}$$

# Example: State estimate (3D-Var)



$$a = x - 2y + 3z$$

$$\begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 1 & -2 & 3 \\ 4 & 0 & -5 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix}$$

Here's the other equation...  $b = 4x - 5z$

Diagram illustrating the forward and adjoint models with matrix representations. The **Forward model** (blue arrow) is represented by the matrix equation:

$$\begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 1 & -2 & 3 \\ 4 & 0 & -5 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix}$$

The **Adjoint model** (purple arrow) is represented by the matrix equation:

$$\begin{pmatrix} \hat{x} \\ \hat{y} \\ \hat{z} \end{pmatrix} = \begin{pmatrix} 1 & 4 \\ -2 & 0 \\ 3 & -5 \end{pmatrix} \begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix}$$

Diagram illustrating the adjoint model with matrix representations. The adjoint model is represented by the matrix equation:

$$\begin{pmatrix} \hat{x} \\ \hat{y} \\ \hat{z} \end{pmatrix} = \begin{pmatrix} 1 & 4 \\ -2 & 0 \\ 3 & -5 \end{pmatrix} \begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix}$$

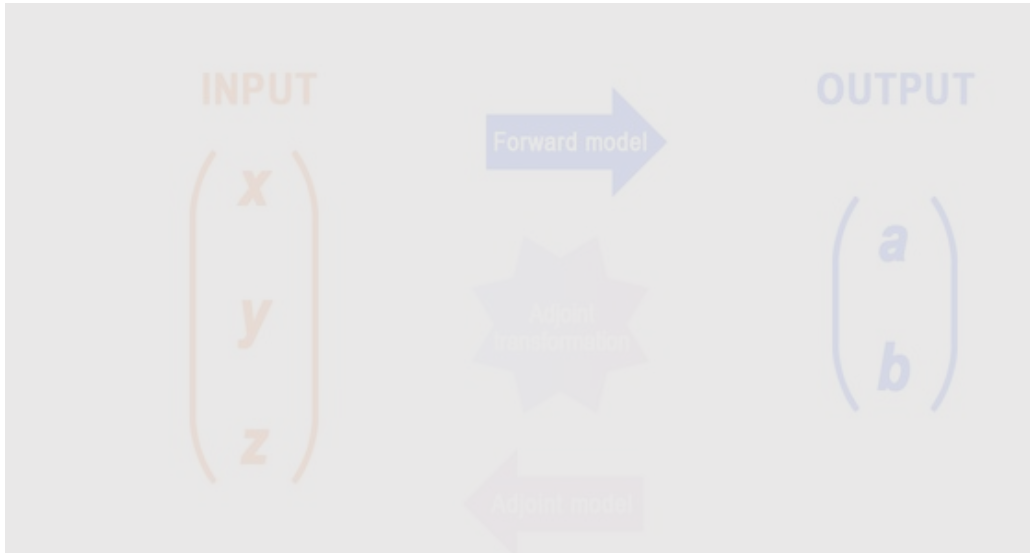
The corresponding state estimates are:

$$\hat{x} = \hat{a} + 4\hat{b}$$

$$\hat{y} = -2\hat{a}$$

$$\hat{z} = 3\hat{a} - 5\hat{b}$$

# Example: State estimate (3D-Var)



$$a = x - 2y + 3z$$

$$\begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 1 & -2 & 3 \\ 4 & 0 & -5 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix}$$

Here's the other equation...  $b = 4x - 5z$

**Forward model**  $\rightarrow$

$$\begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 1 & -2 & 3 \\ 4 & 0 & -5 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix}$$

**Adjoint model**  $\leftarrow$

$$\begin{pmatrix} \hat{x} \\ \hat{y} \\ \hat{z} \end{pmatrix} = \begin{pmatrix} 1 & 4 \\ -2 & 0 \\ 3 & -5 \end{pmatrix} \begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix}$$

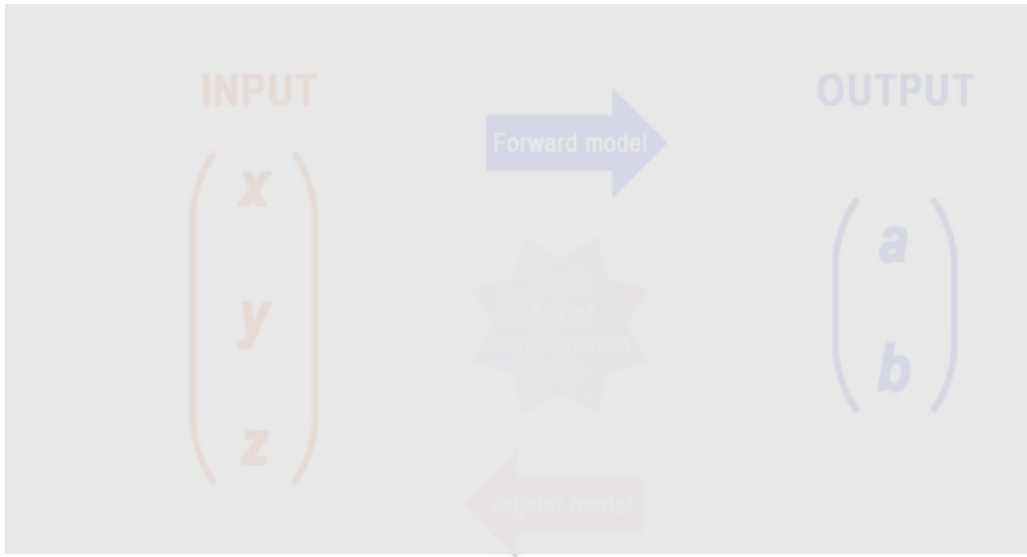
$$\begin{pmatrix} \hat{x} \\ \hat{y} \\ \hat{z} \end{pmatrix} = \begin{pmatrix} 1 & 4 \\ -2 & 0 \\ 3 & -5 \end{pmatrix} \begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix}$$

$$\hat{x} = \hat{a} + 4\hat{b}$$

$$\hat{y} = -2\hat{a}$$

$$\hat{z} = 3\hat{a} - 5\hat{b}$$

# Example: State estimate (3D-Var)



$$a = x - 2y + 3z$$

$$\begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 1 & -2 & 3 \\ 4 & 0 & -5 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix}$$

Here's the other equation...  $b = 4x - 5z$

Forward model:  $\begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 1 & -2 & 3 \\ 4 & 0 & -5 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix}$

Adjoint model:  $\begin{pmatrix} \hat{x} \\ \hat{y} \\ \hat{z} \end{pmatrix} = \begin{pmatrix} 1 & 4 \\ -2 & 0 \\ 3 & -5 \end{pmatrix} \begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix}$

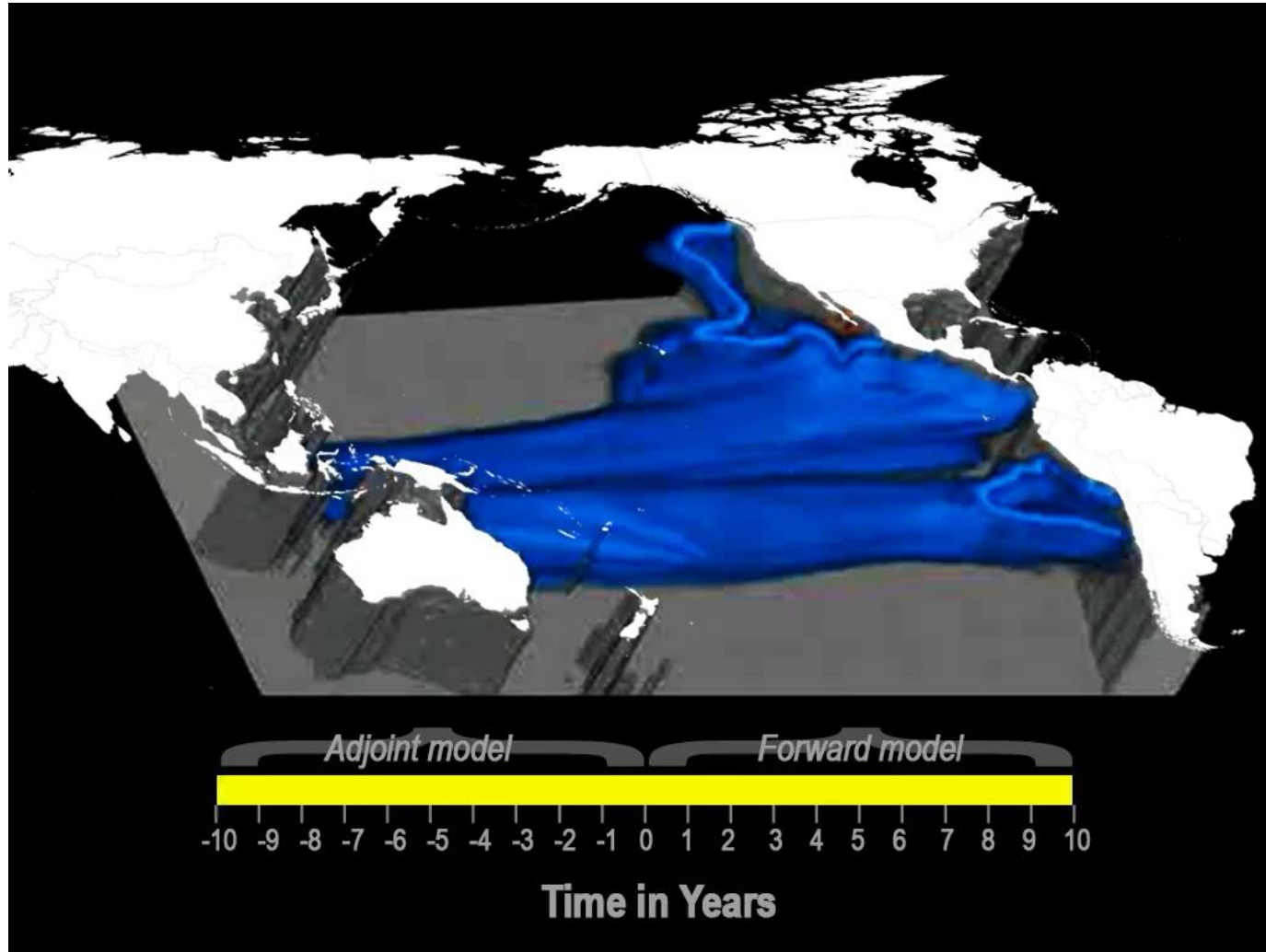
$$\begin{pmatrix} \hat{x} \\ \hat{y} \\ \hat{z} \end{pmatrix} = \begin{pmatrix} 1 & 4 \\ -2 & 0 \\ 3 & -5 \end{pmatrix} \begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix}$$

$$\hat{x} = \hat{a} + 4\hat{b}$$

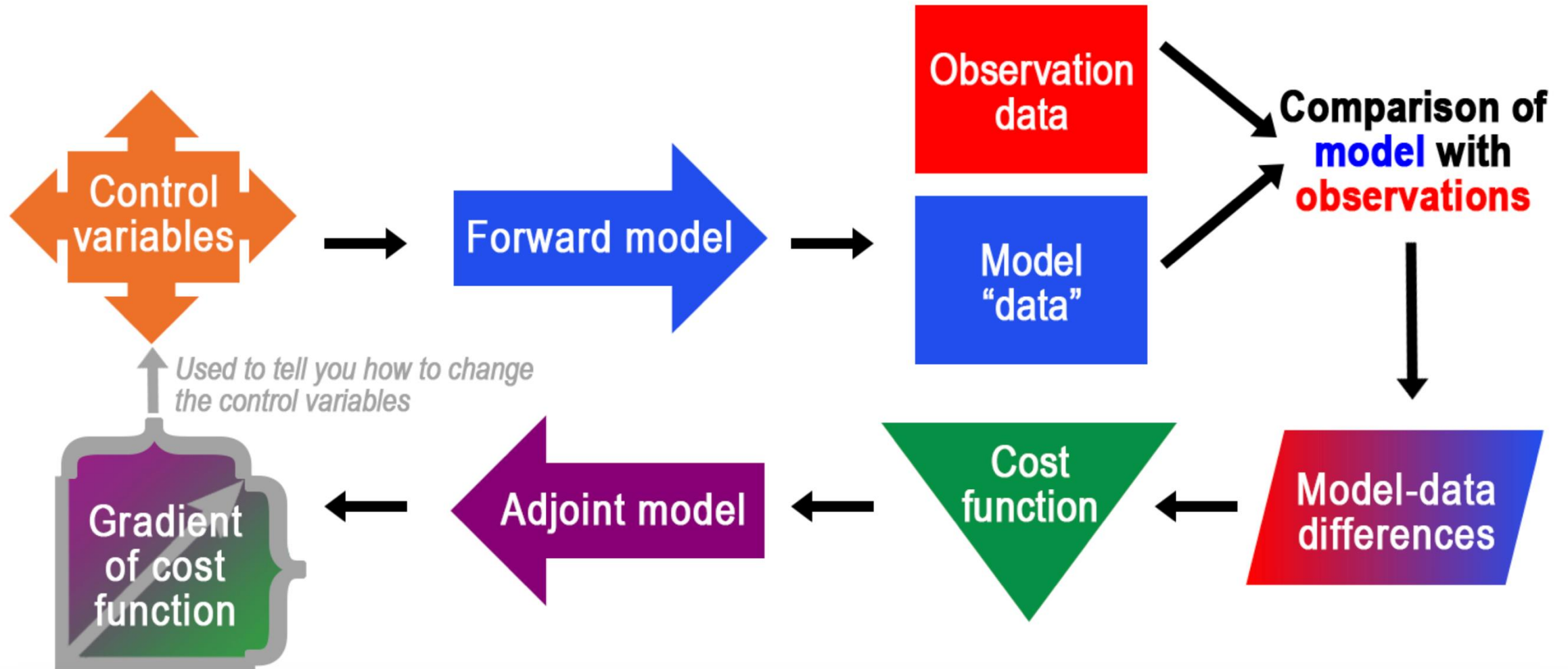
$$\hat{y} = -2\hat{a}$$

$$\hat{z} = 3\hat{a} - 5\hat{b}$$

# Example: State estimate (3D-Var)



# Example: State estimate (3D-Var)



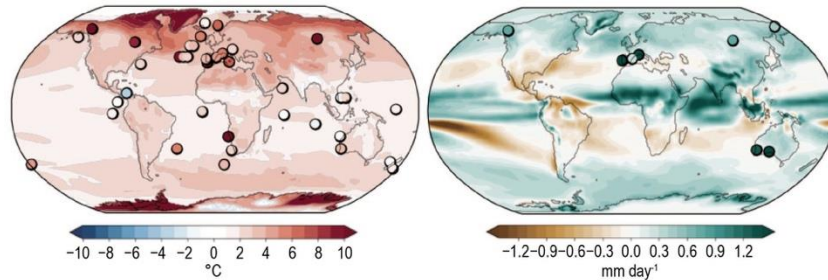
# Example: multiple data sources

## Paleoclimate Models

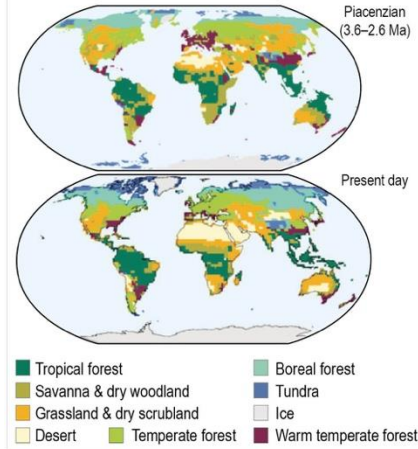
## Proxies and Archives

Climate indicators of the mid-Pliocene Warm Period

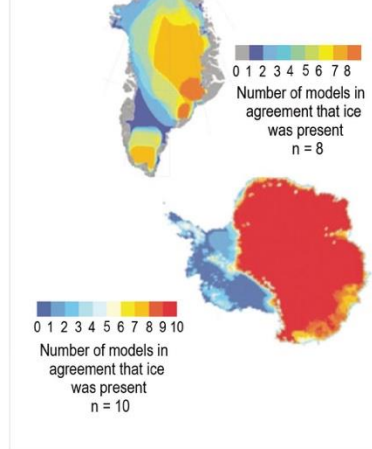
(a) Surface air temperature and precipitation rate anomalies relative to 1850–1900



(b) Changes in vegetation from the Piacenzian to present day



(c)



"Travertine speleothem (Crystal Cave, Main Island, Bermuda) 1" by James S. John is licensed under [CC-BY 2.0](#)



"Tree rings" by Out of the Fire Ring is licensed under [CC-BY 2.0](#)



"A volcanic ash layer in the WAIS Divide ice core. Volcanic markers like these were used in the new study to synchronize ice cores from across Antarctica." by Oregon State University is licensed under [CC-BY-SA 2.0](#)



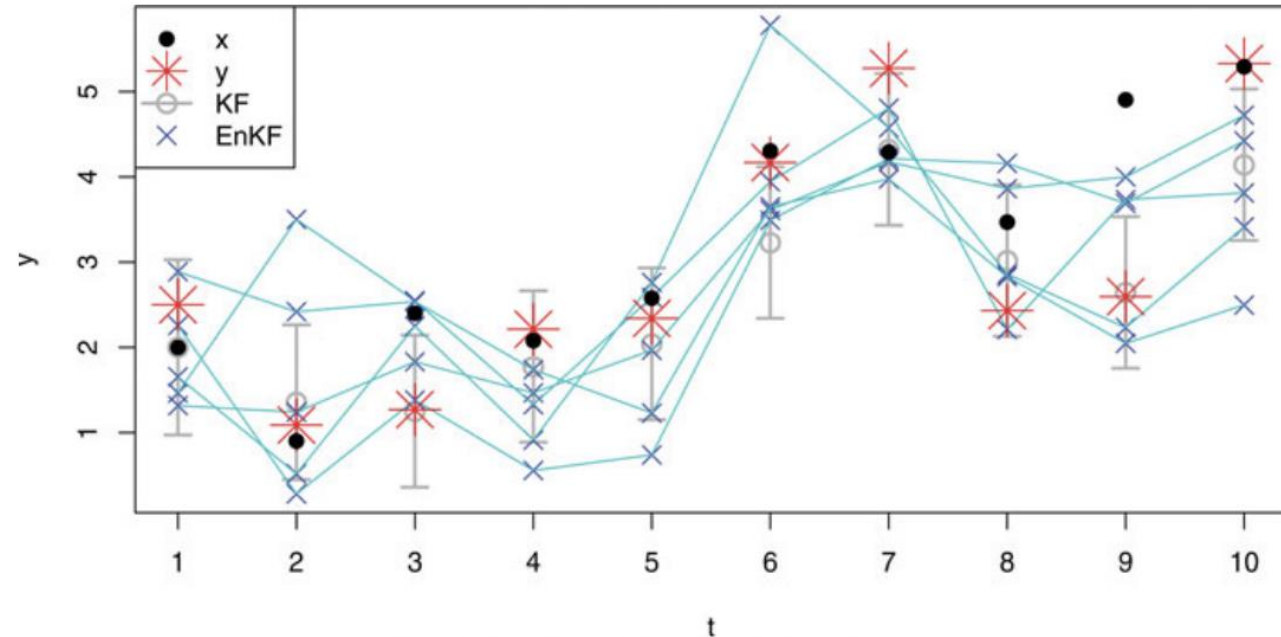
"2016\_Lake sediment core, Forillon Lakes, Gifford Pinchot National Forest, Washington" by USDA Forest Service is marked with [Public Domain Mark 1.0](#)



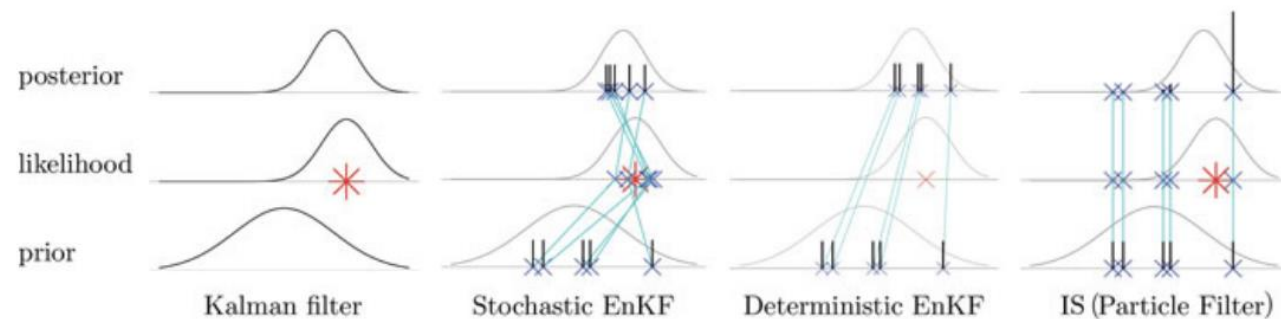
"Diploria fossil brain coral on Devil's Point Hardground (Cockburn Town Member, Grotto Beach Formation, Upper Pleistocene, ~120-123 ka, Cockburn Town Fossil Reef, San Salvador Island, Bahamas) 3" by James S. John is licensed under [CC-BY 2.0](#)

# Example: multiple data sources

## Ensemble Kalman Filter



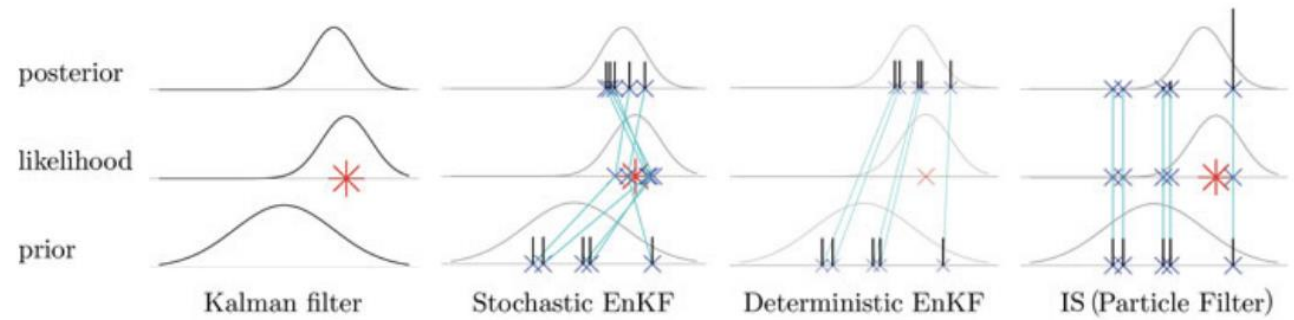
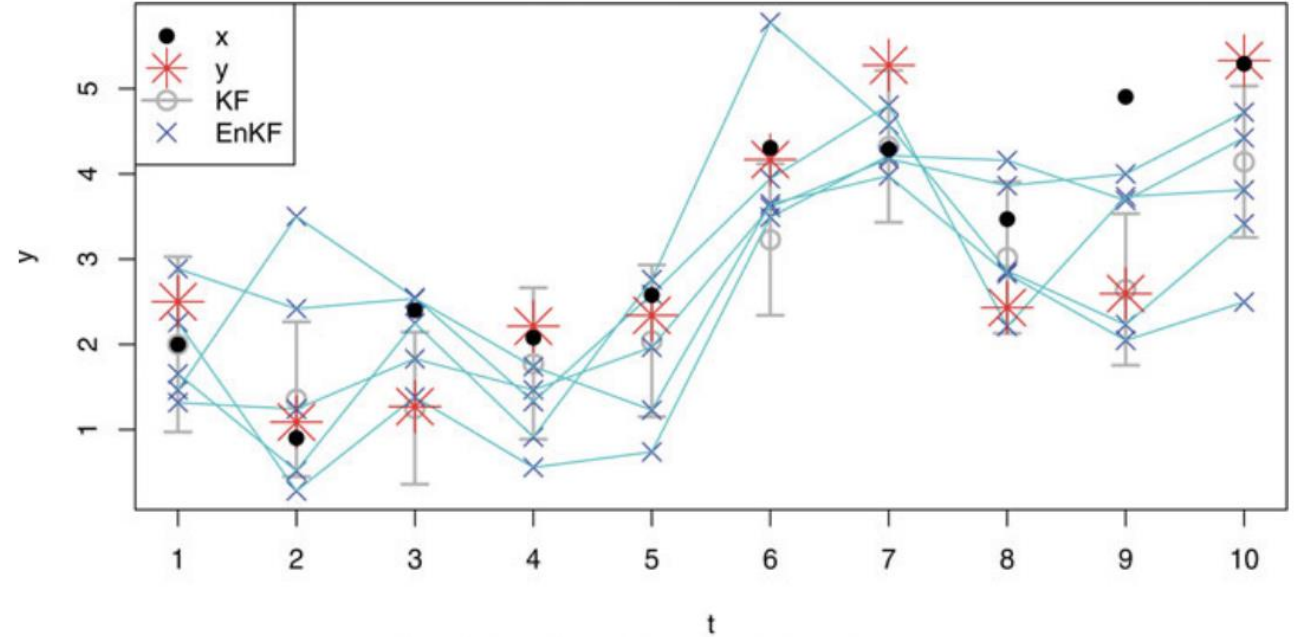
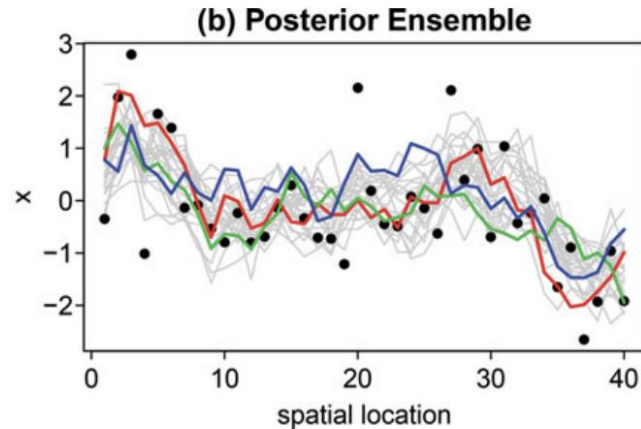
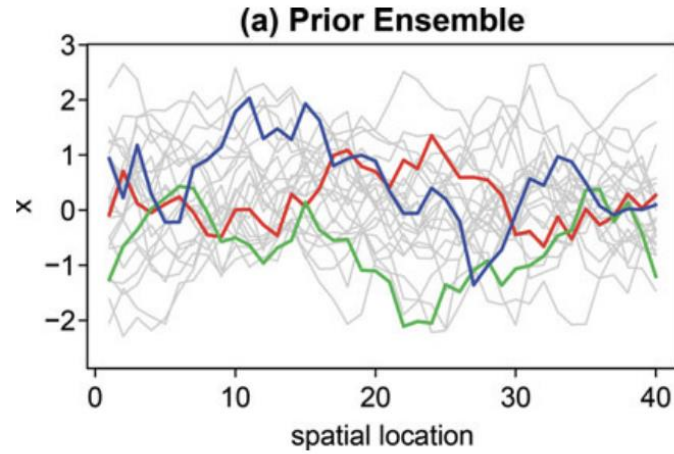
(a) State, observation, Kalman filter, and EnKF for 10 time points



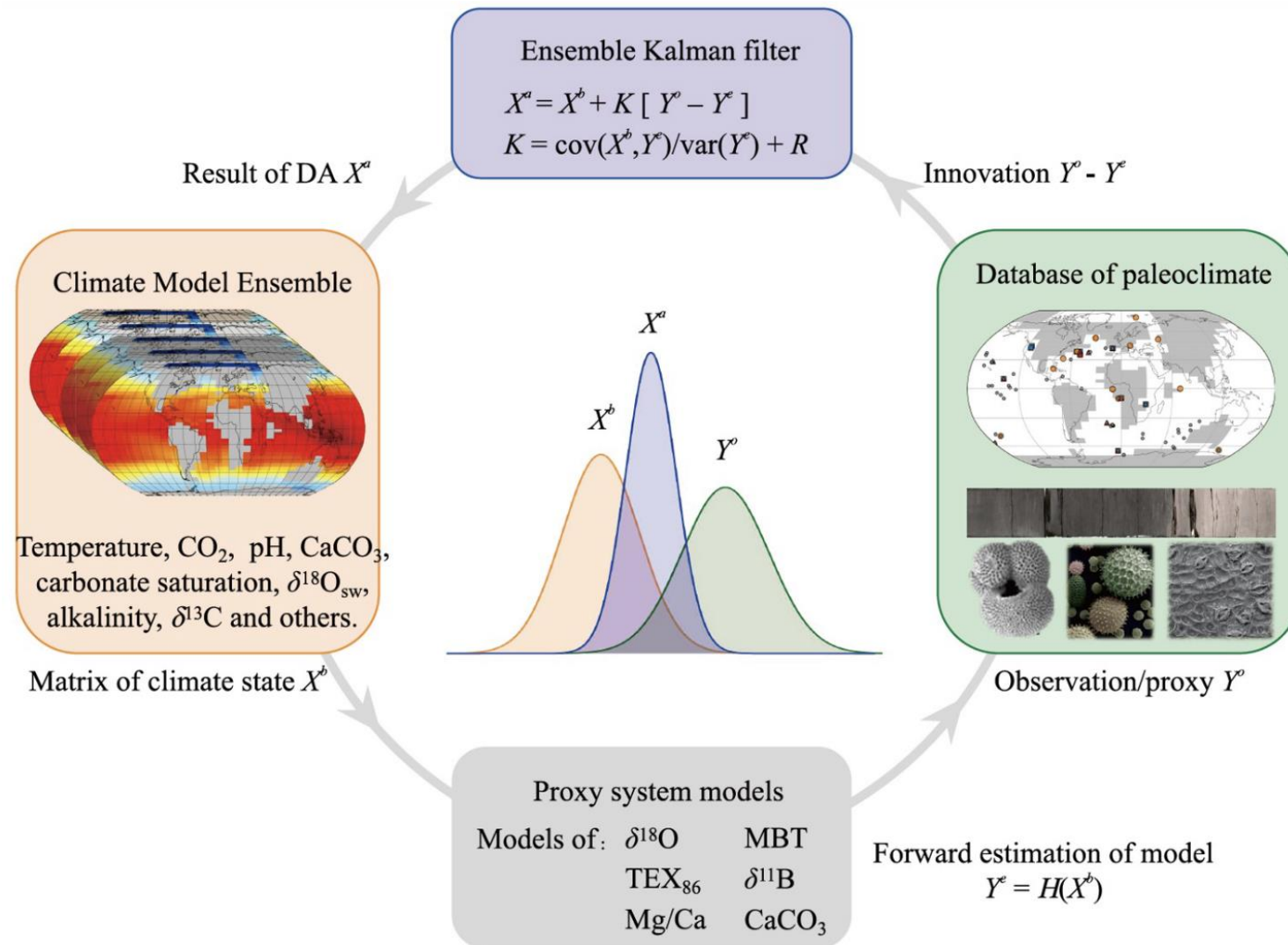
(b) Illustration of different updating schemes at time point  $t = 1$

# Example: multiple data sources

## Ensemble Kalman Filter



# Example: multiple data sources

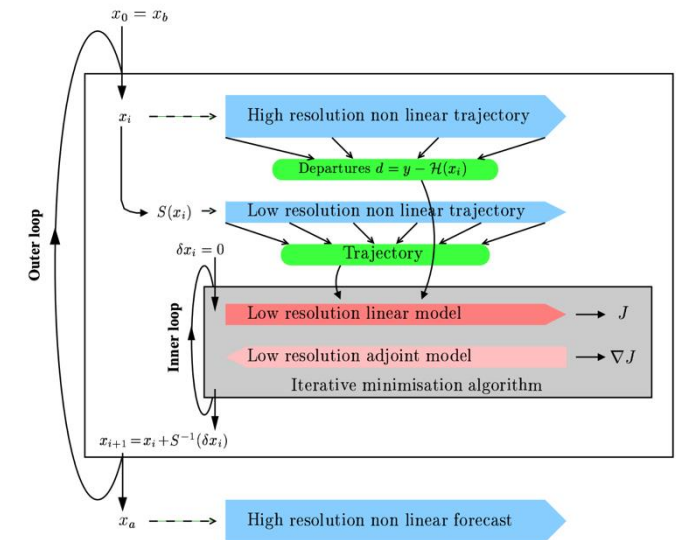


# Data assimilation with AI

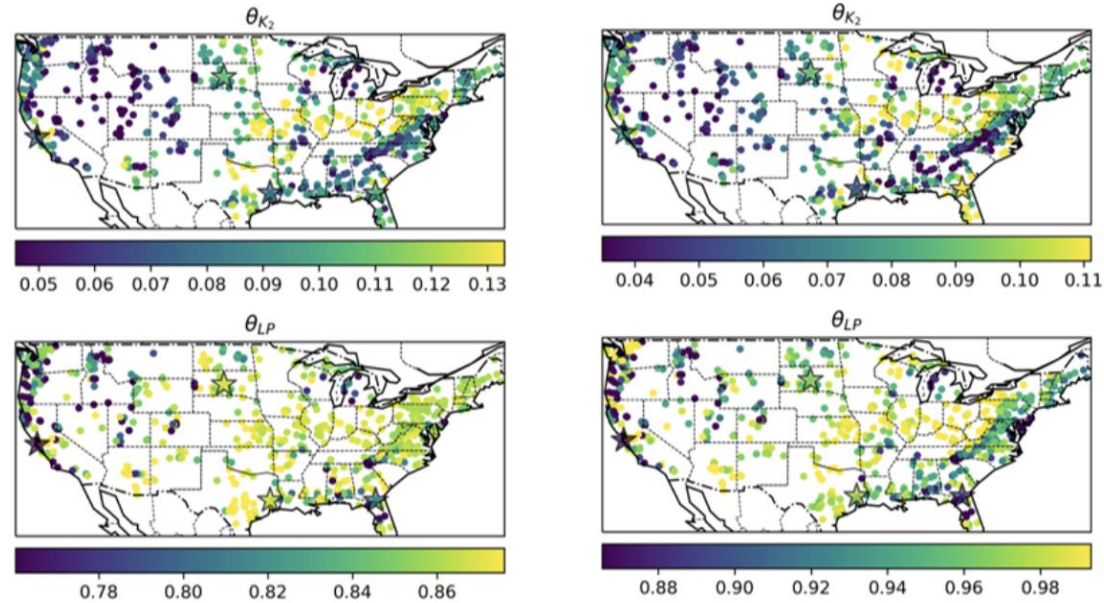
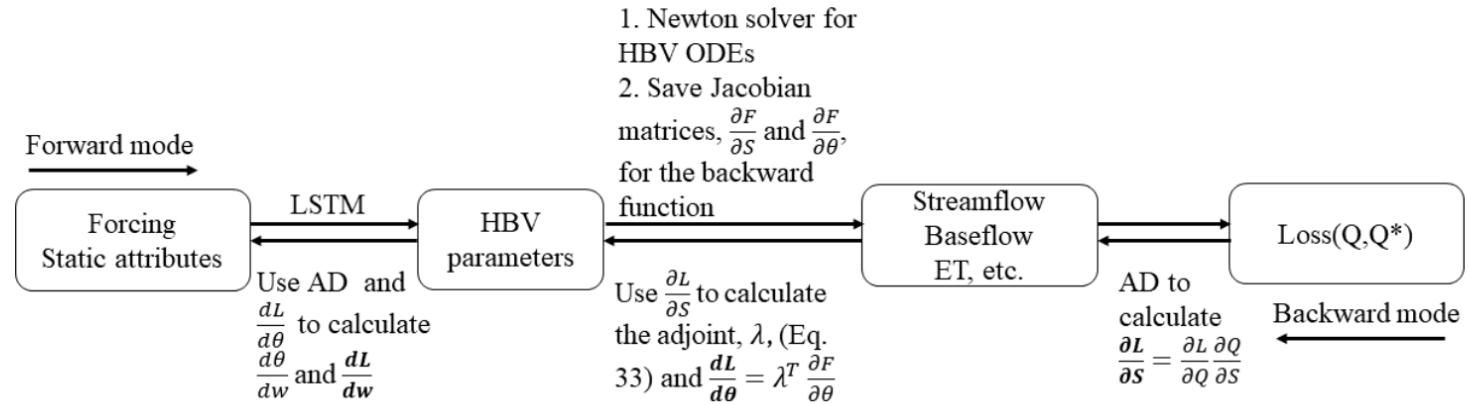
# Data assimilation with AI

Data assimilation is naturally suited for AI integration

- It is already an optimization problem!
- Includes backprop in adjoint methods
- Data-driven task i.e., mostly agnostic to underlying physics
- Potential for big impact in climate: inference is always challenge due to limited sampling



# Example: deep learning with adjoint methods



(a) Sequential model

(b) Adjoint model

# Example: Combined with EnKF

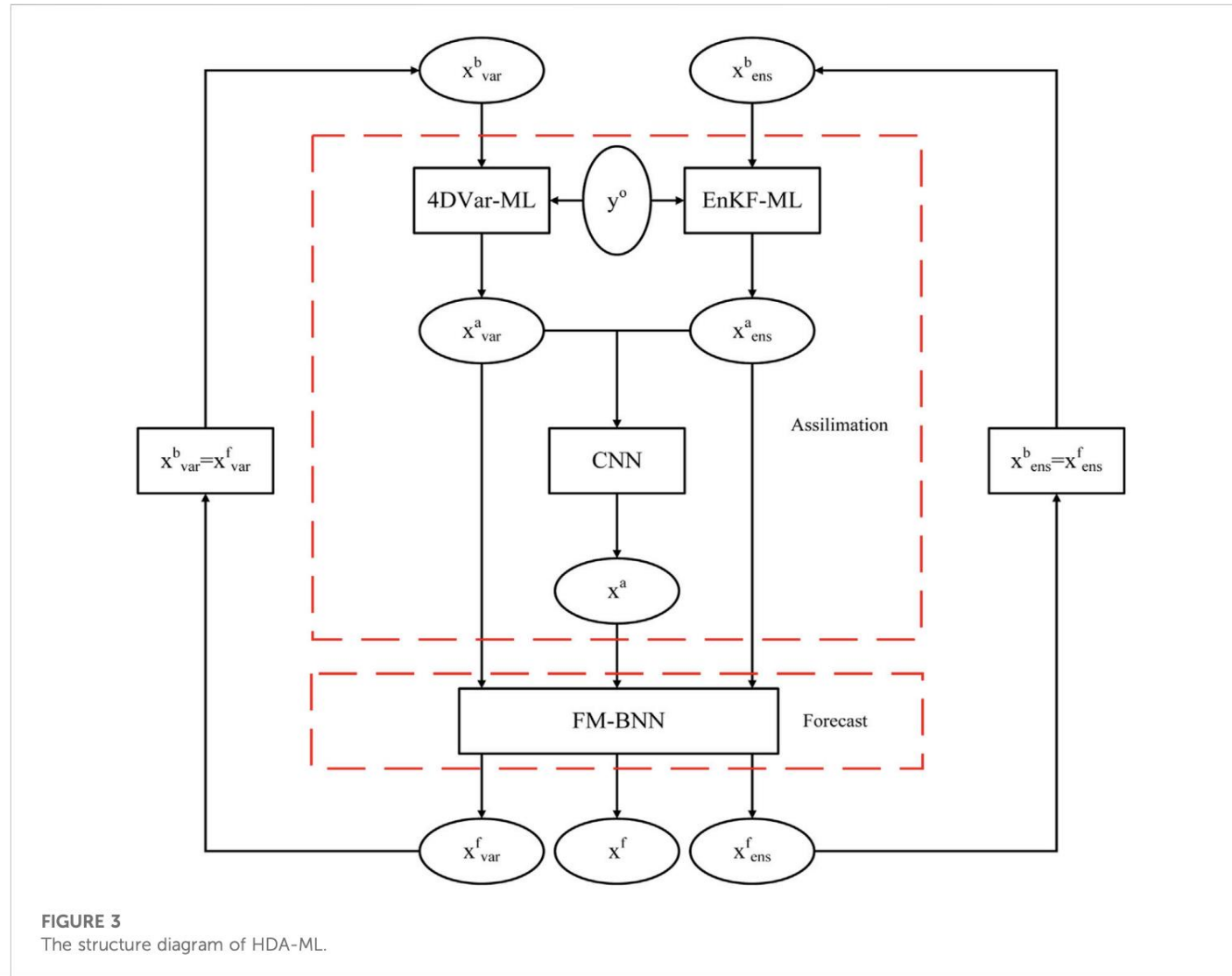


FIGURE 3  
The structure diagram of HDA-ML.

# Example: reduced order modeling for 3D/4D-Var

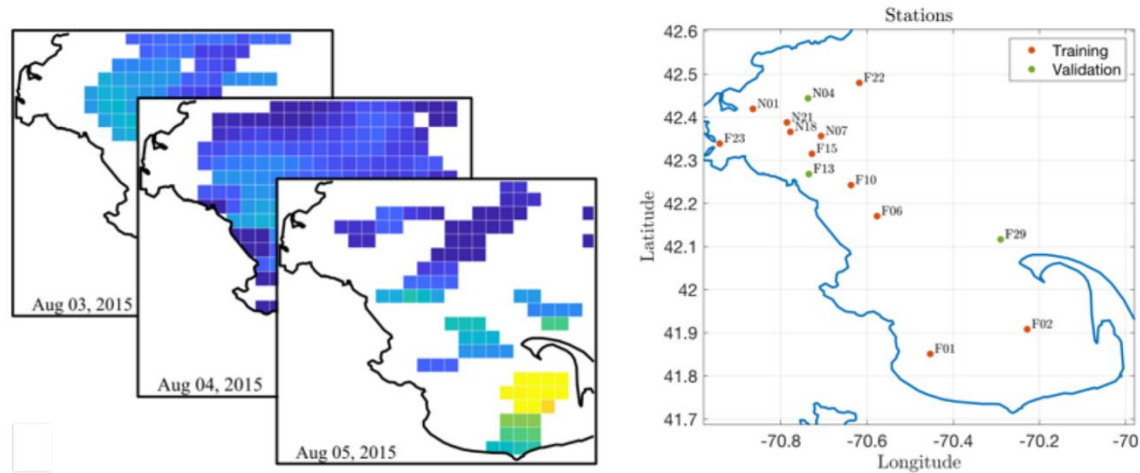


Figure 2: **Sensor Data.** The low fidelity data (satellite (left)) is only available on days with low cloud coverage. The high fidelity data (buoys (right)) is local in space and sparse.

# Example: reduced order modeling for 3D/4D-Var

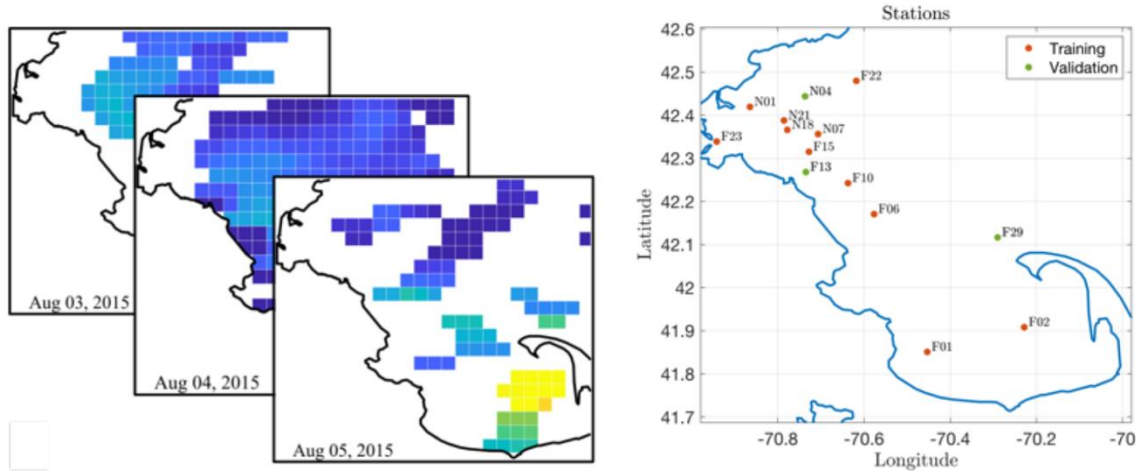
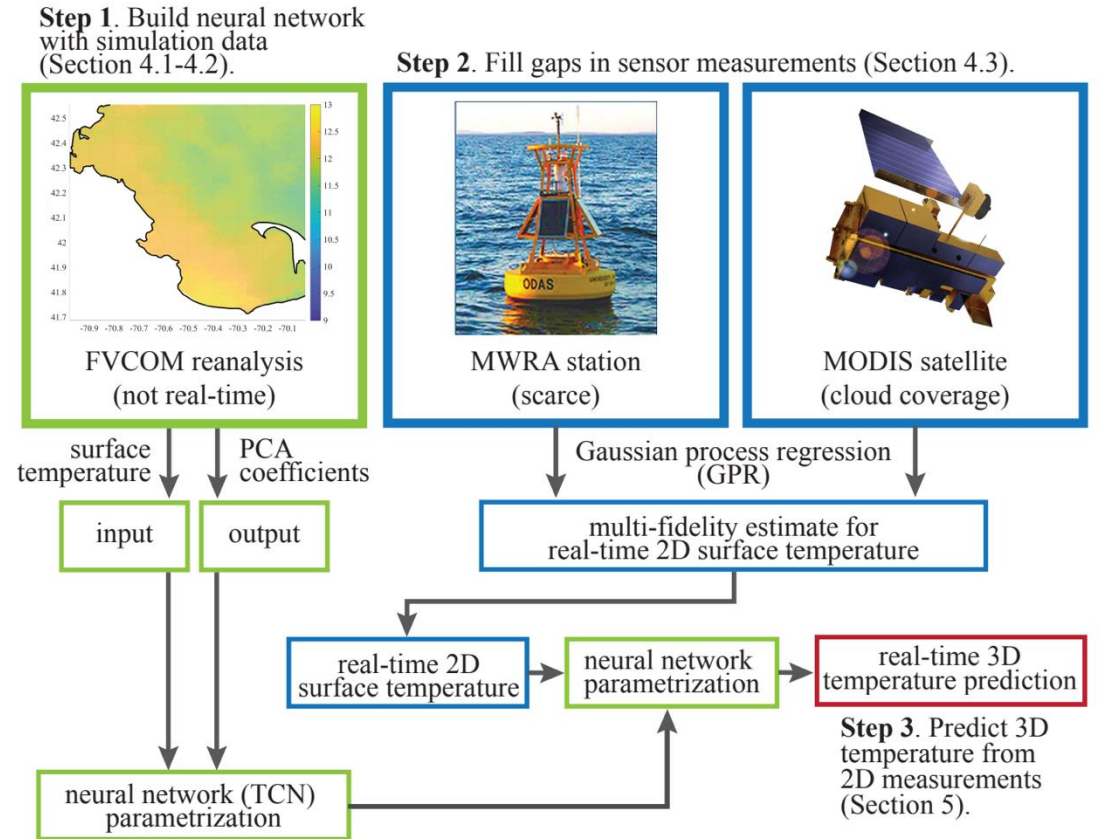
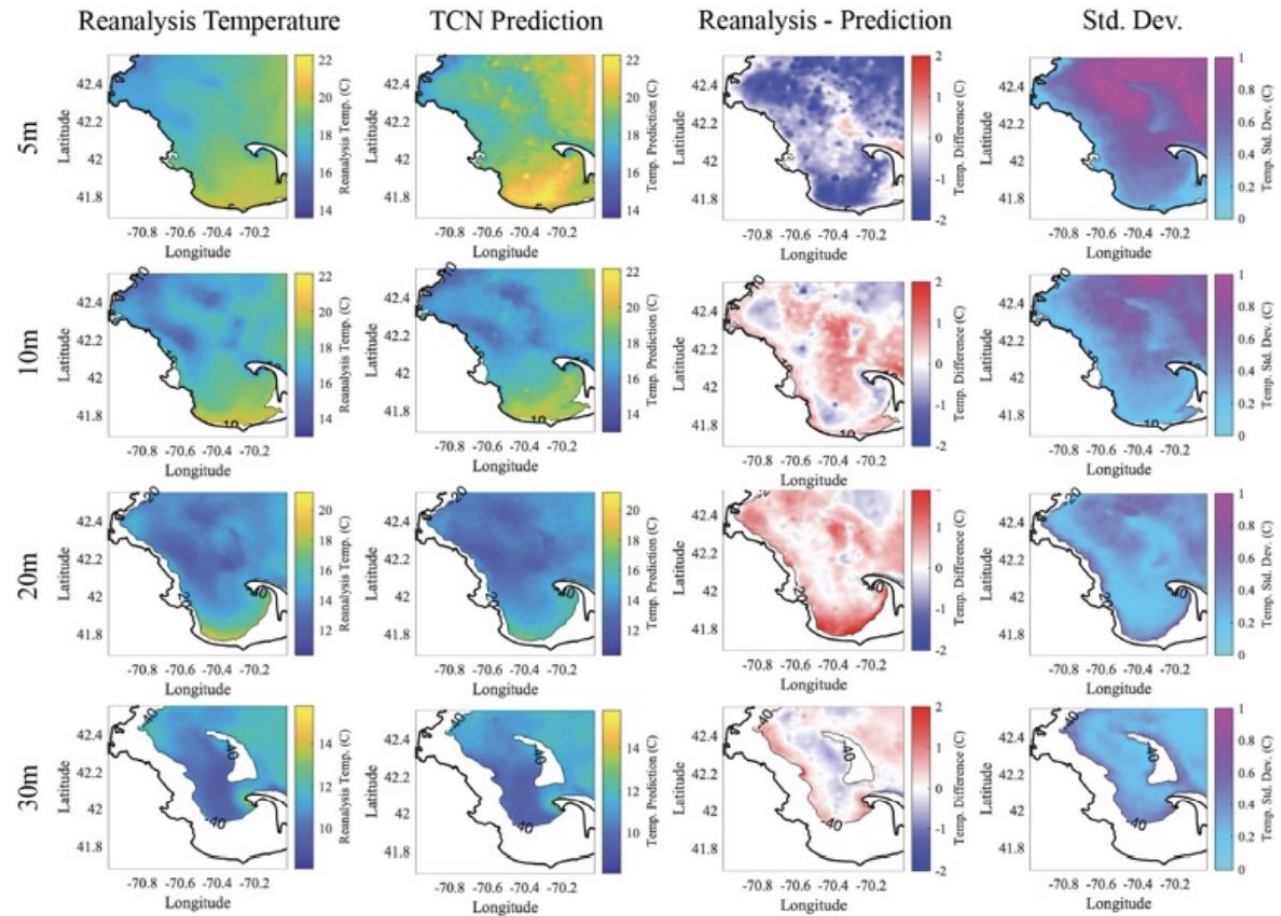


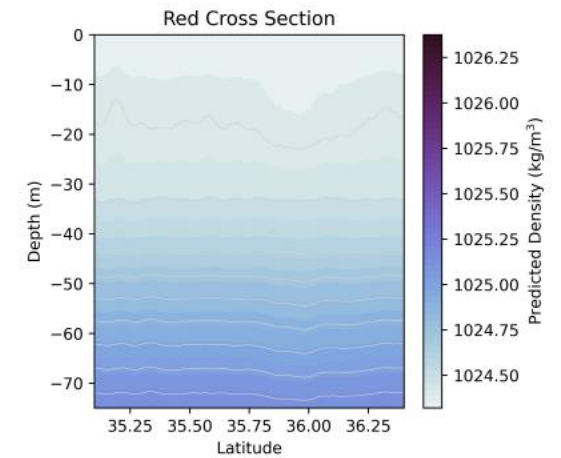
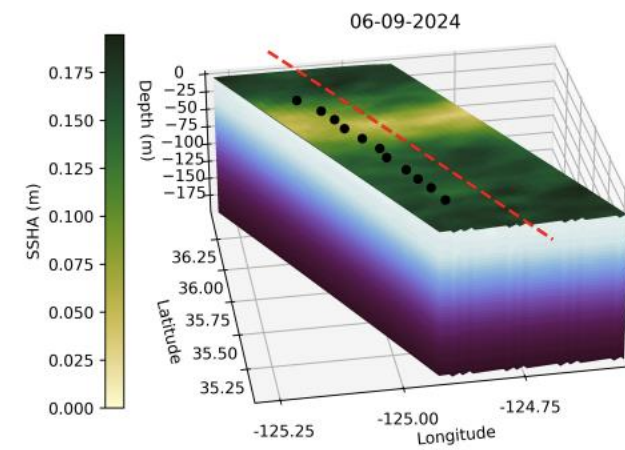
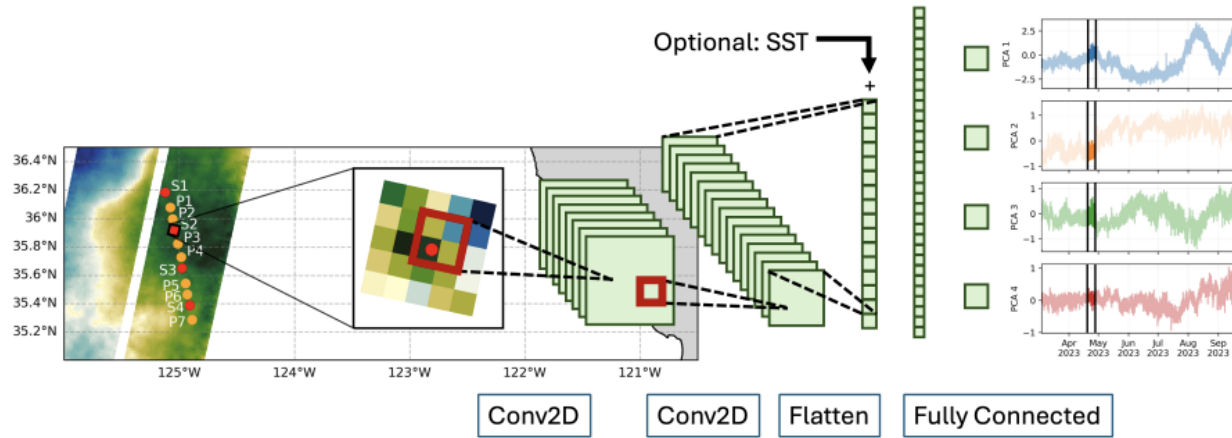
Figure 2: **Sensor Data.** The low fidelity data (satellite (left)) is only available on days with low cloud coverage. The high fidelity data (buoys (right)) is local in space and sparse.



# Example: reduced order modeling for 3D/4D-Var

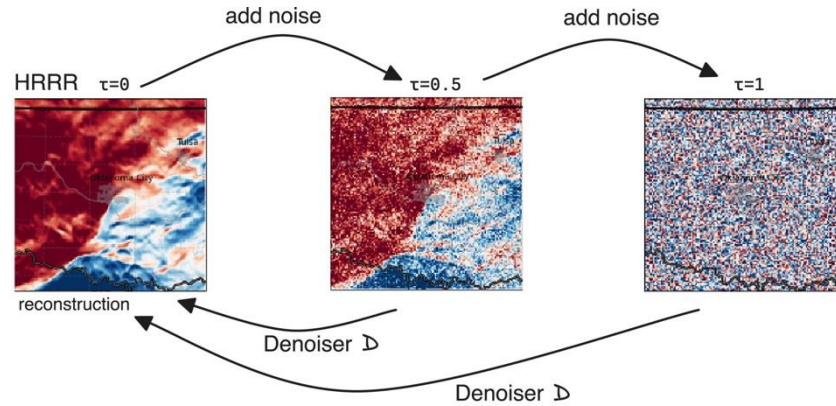


# Example: reduced order modeling for 3D/4D-Var

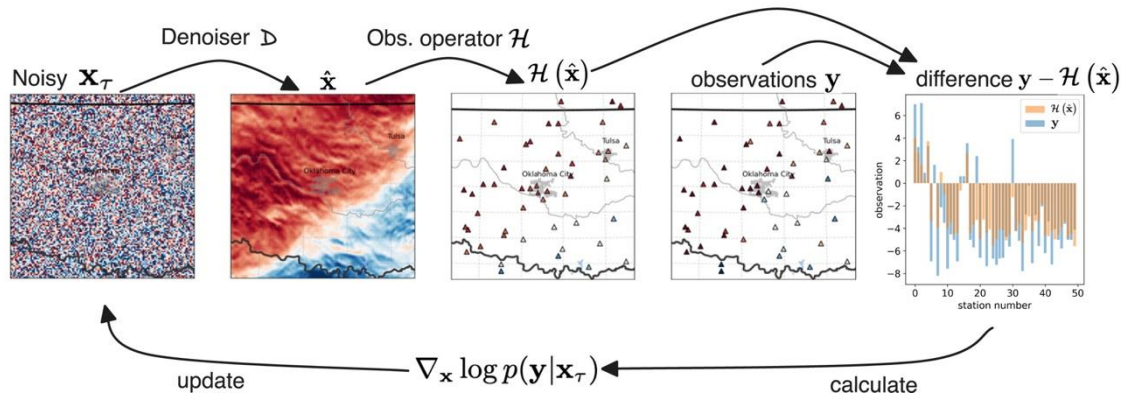


# Example: generative modeling

Denoiser training

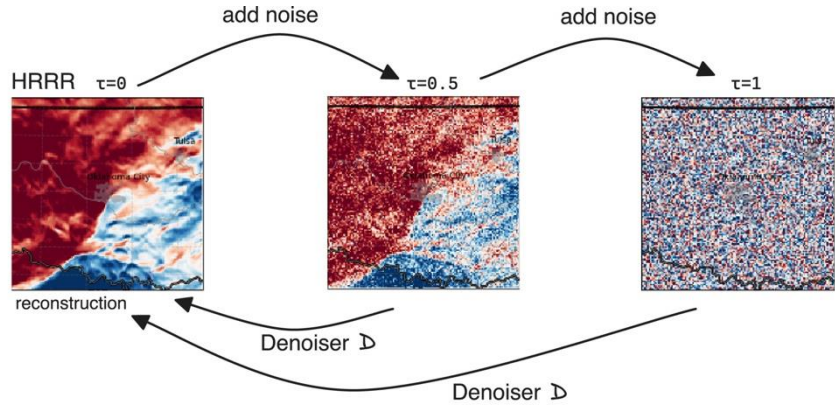


Data assimilation

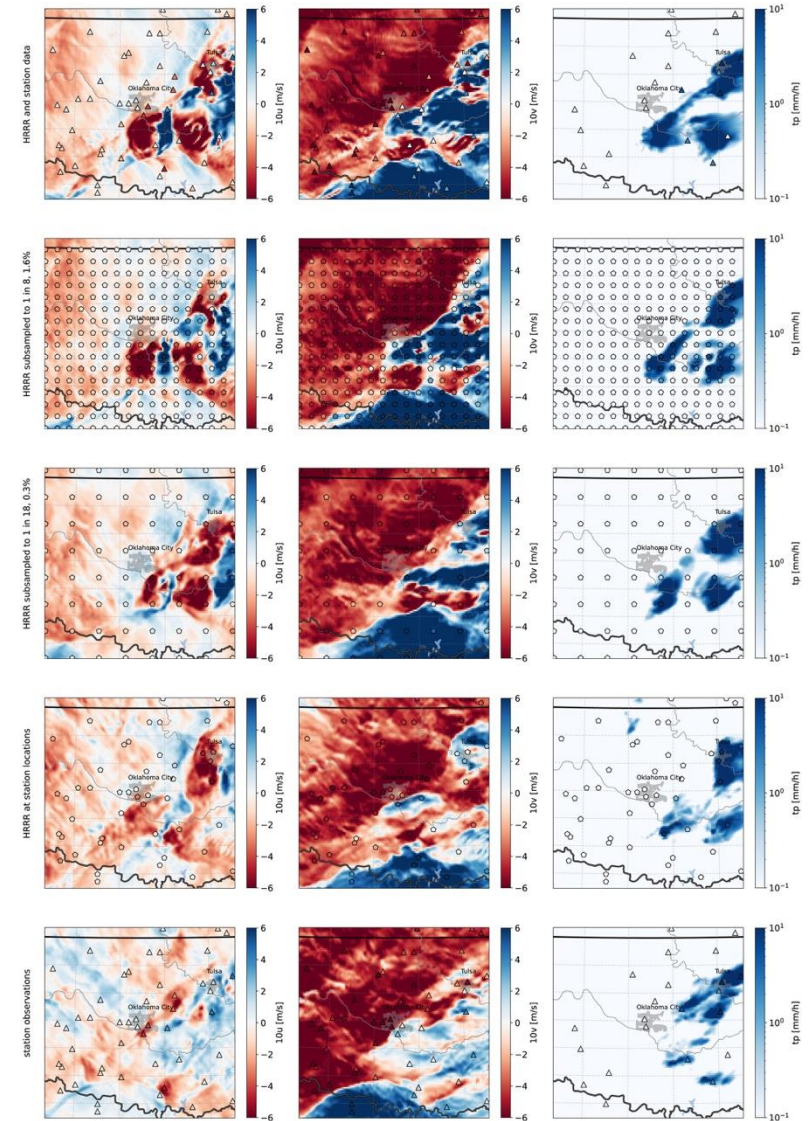
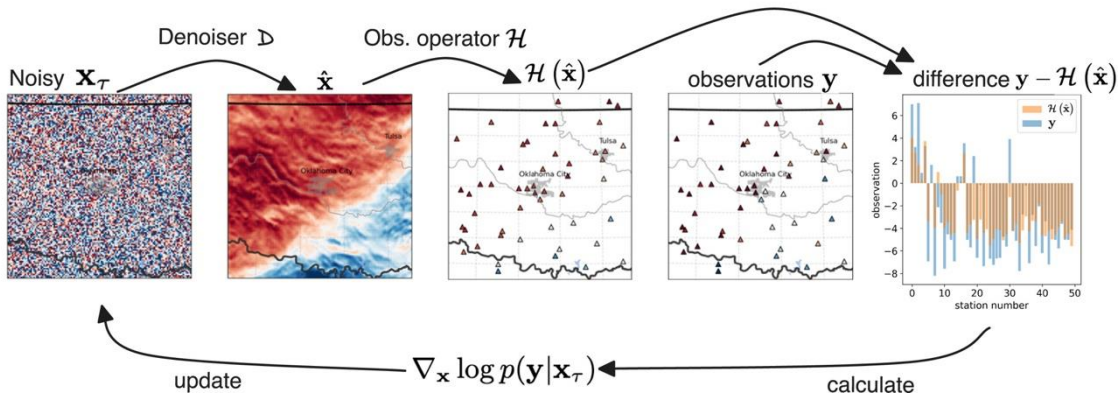


# Example: generative modeling

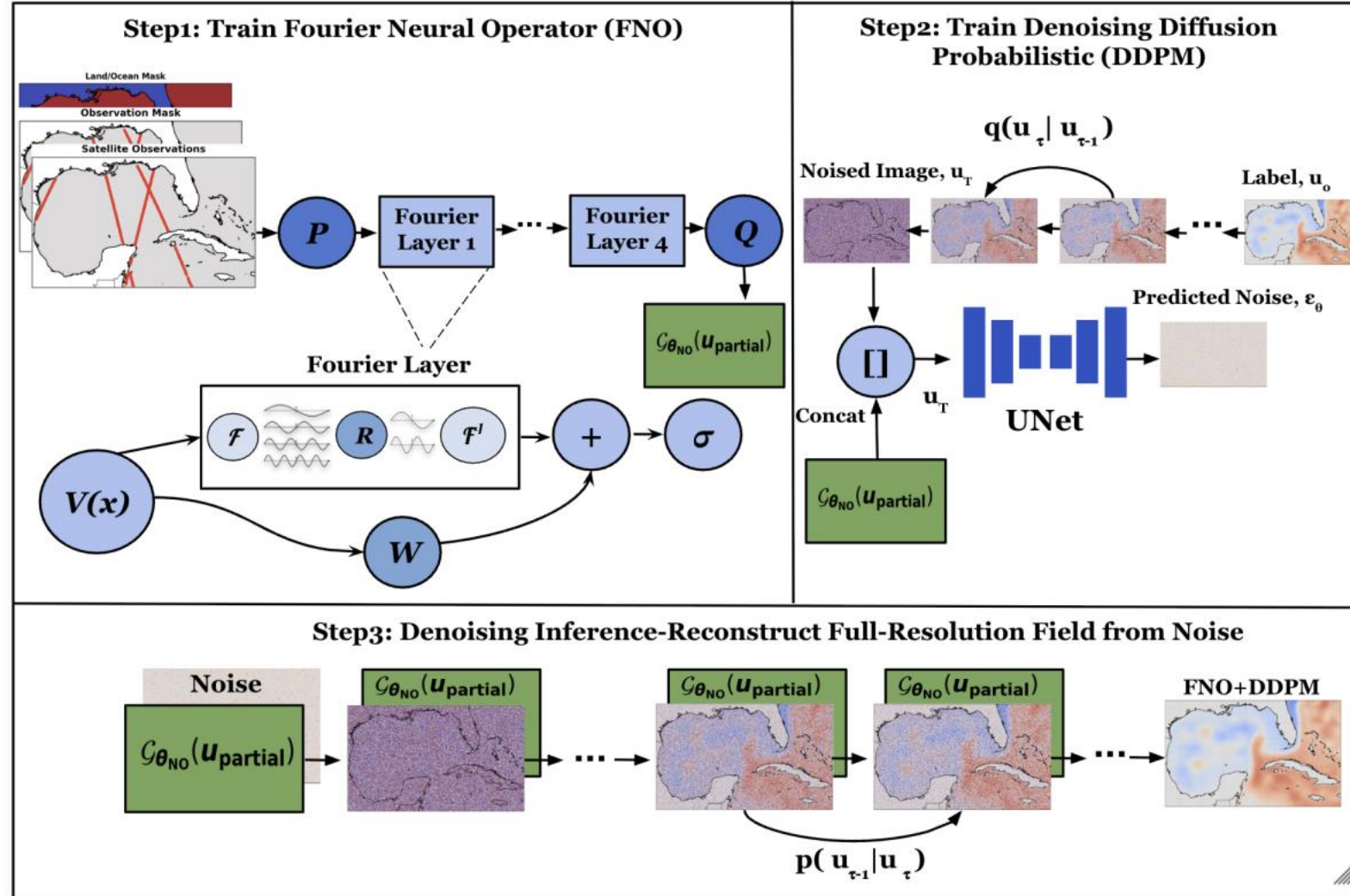
Denoiser training



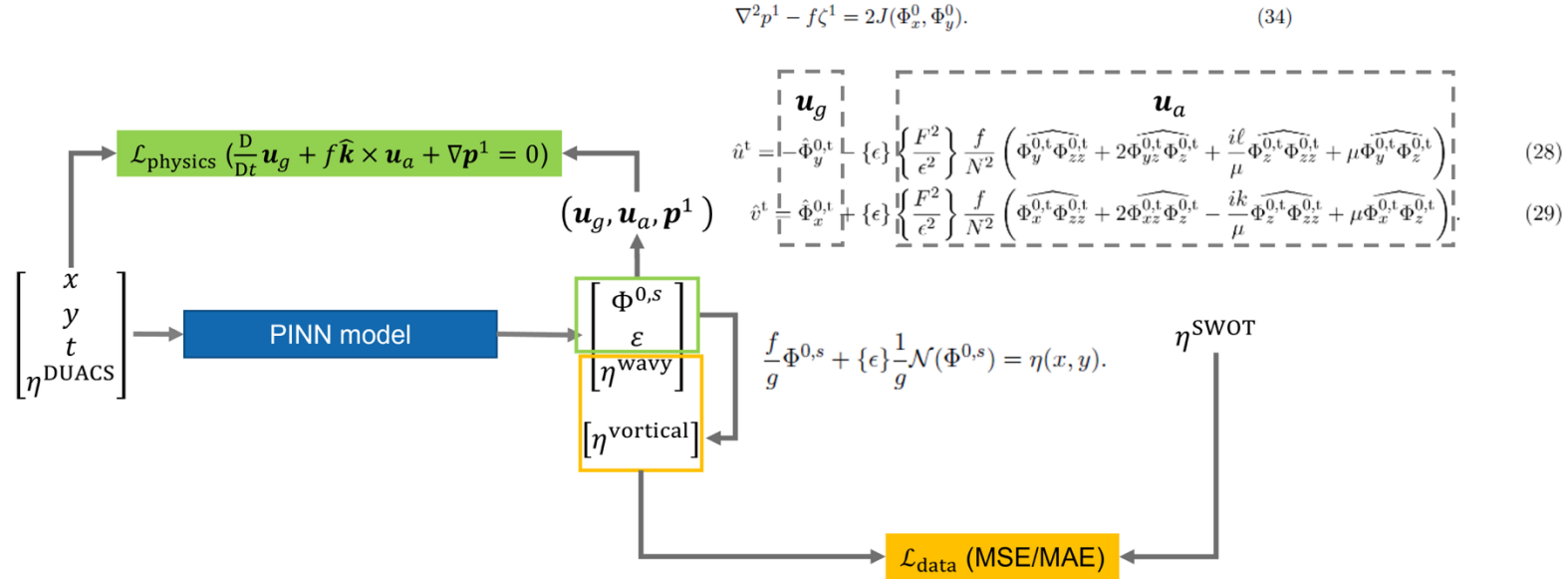
Data assimilation



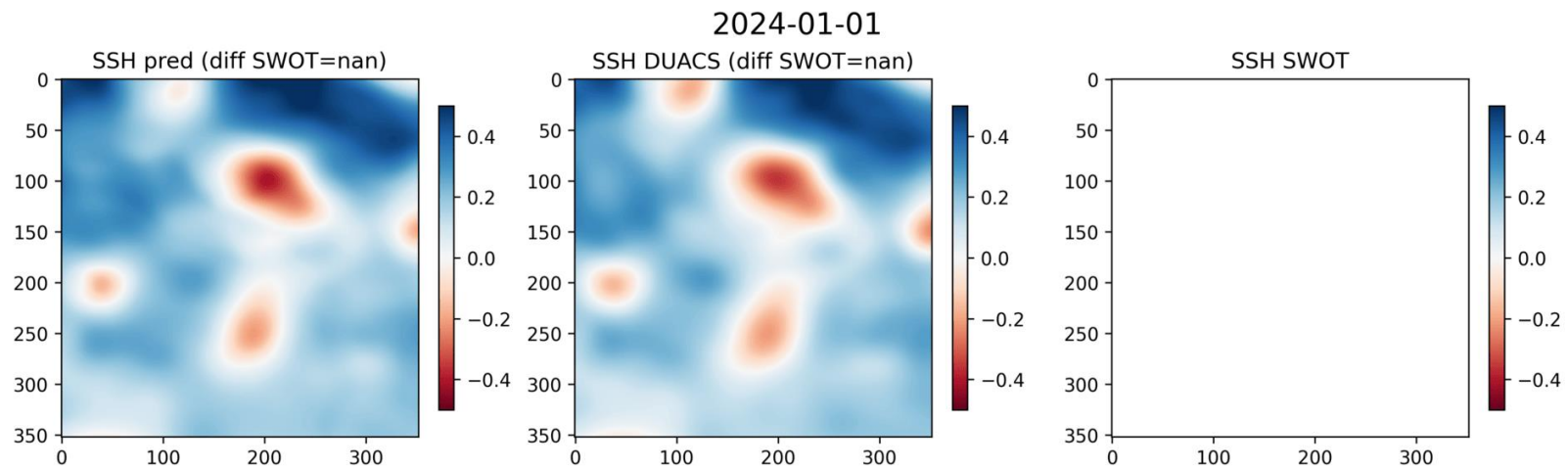
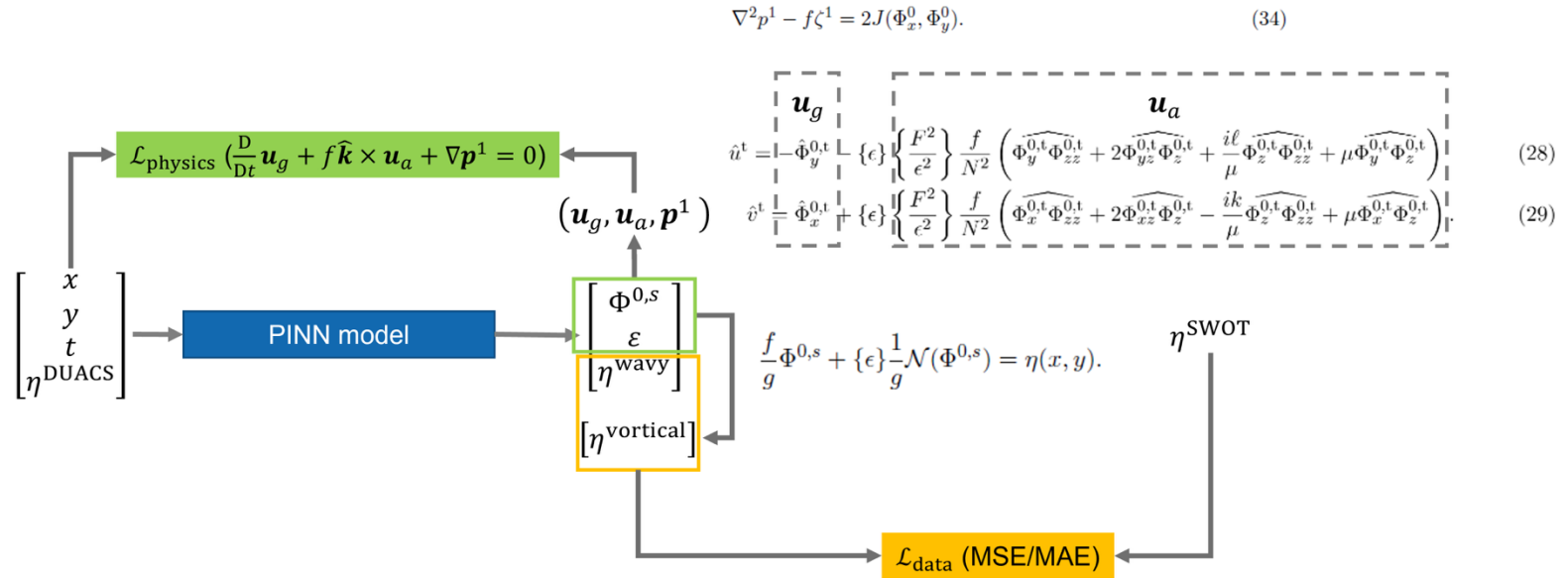
# Example: generative modeling



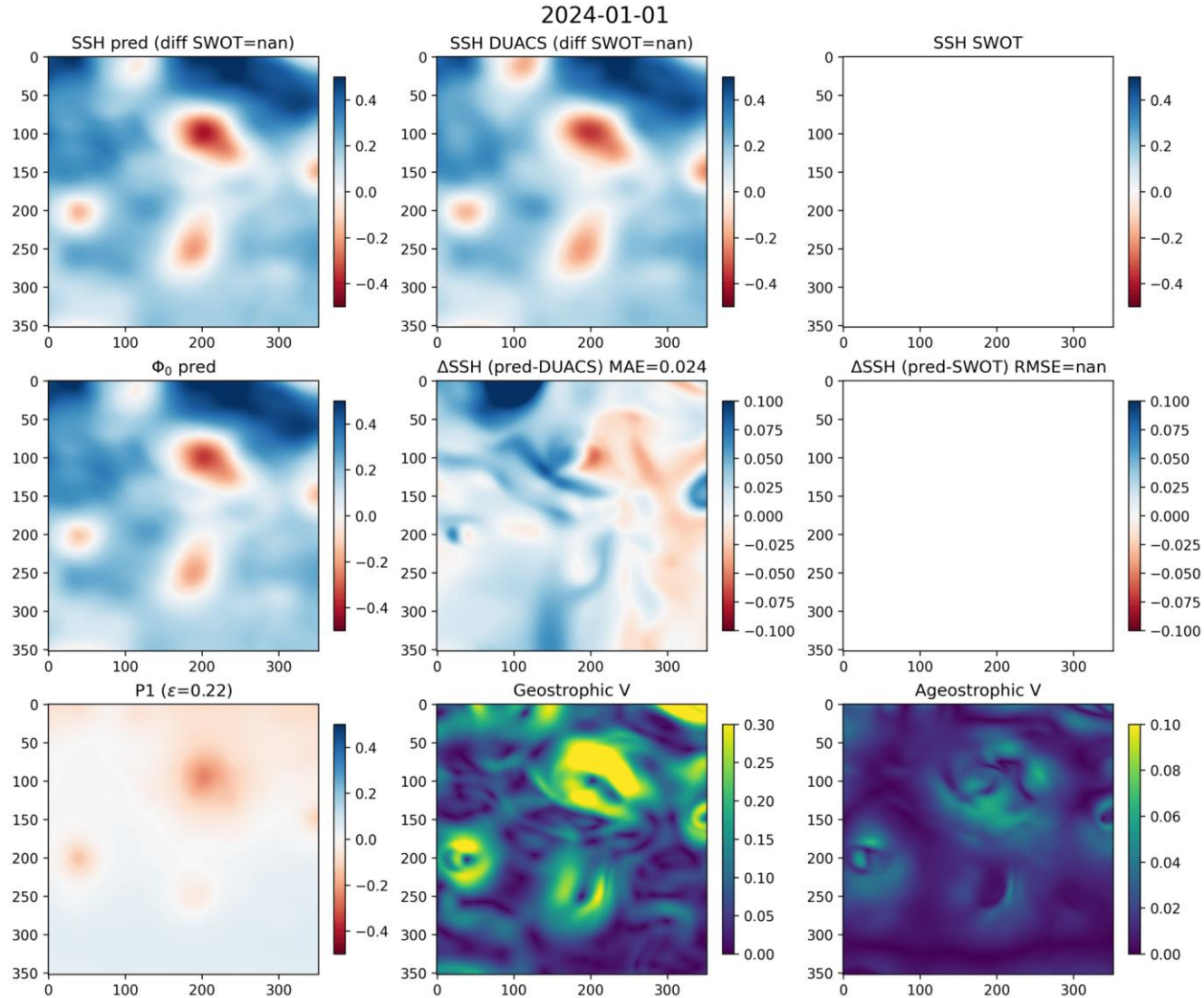
# Example: physics informed neural network (PINN)



# Example: physics informed neural network (PINN)

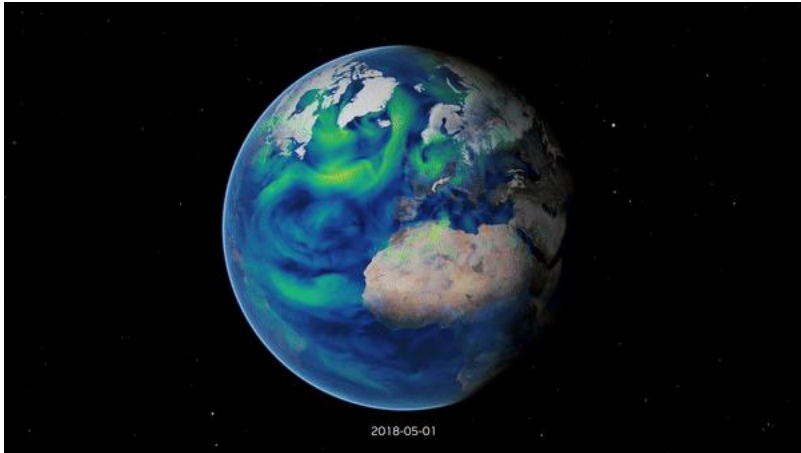


# Example: physics informed neural network (PINN)

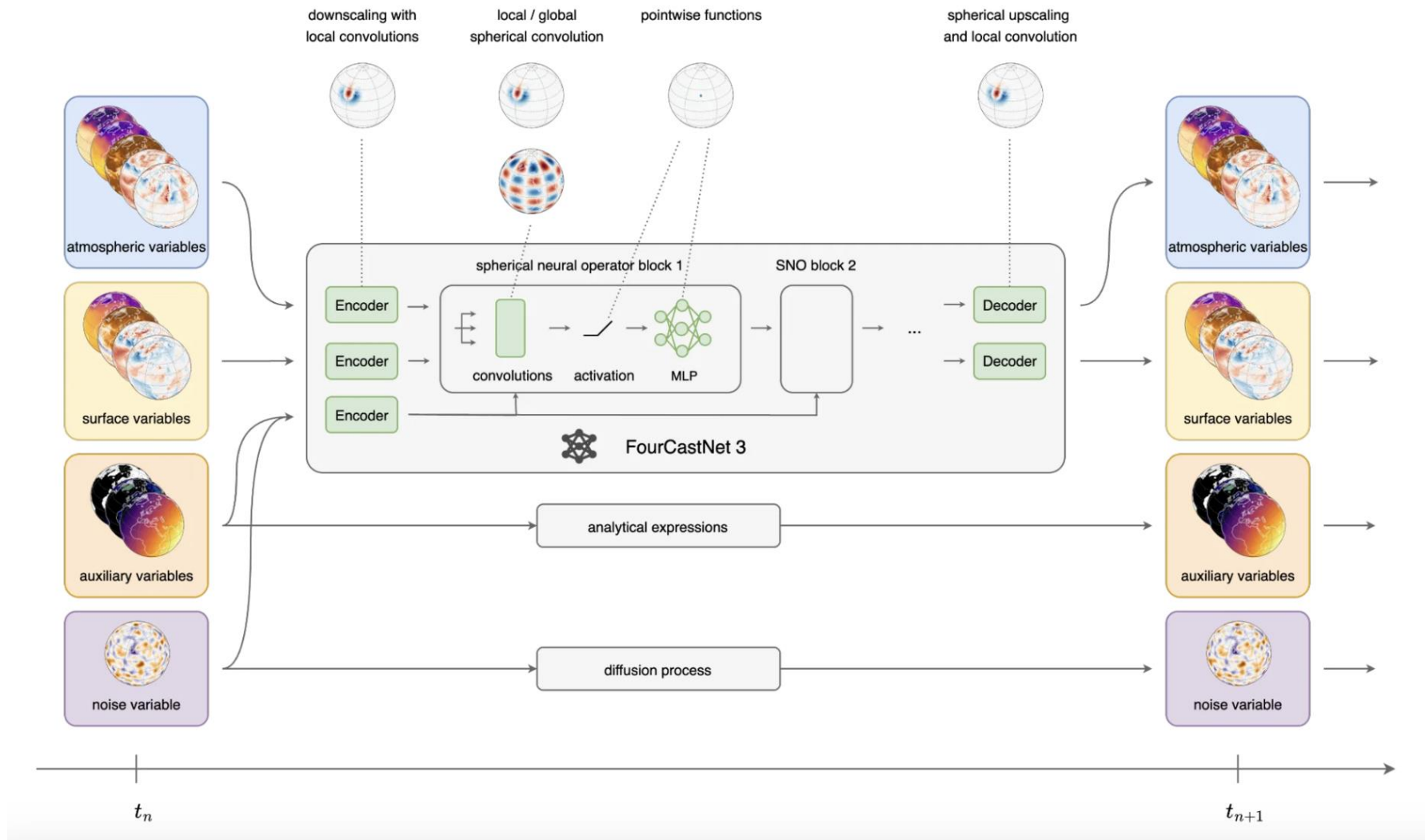


# Example: Forecasting

FourCastNet3

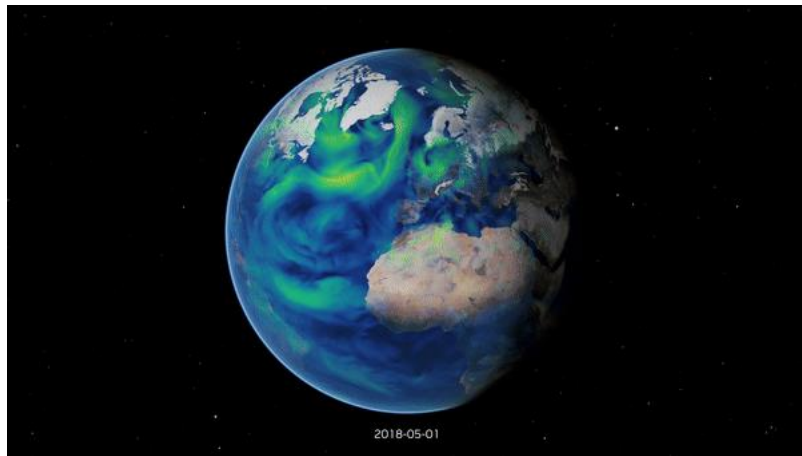


# Example: Forecasting

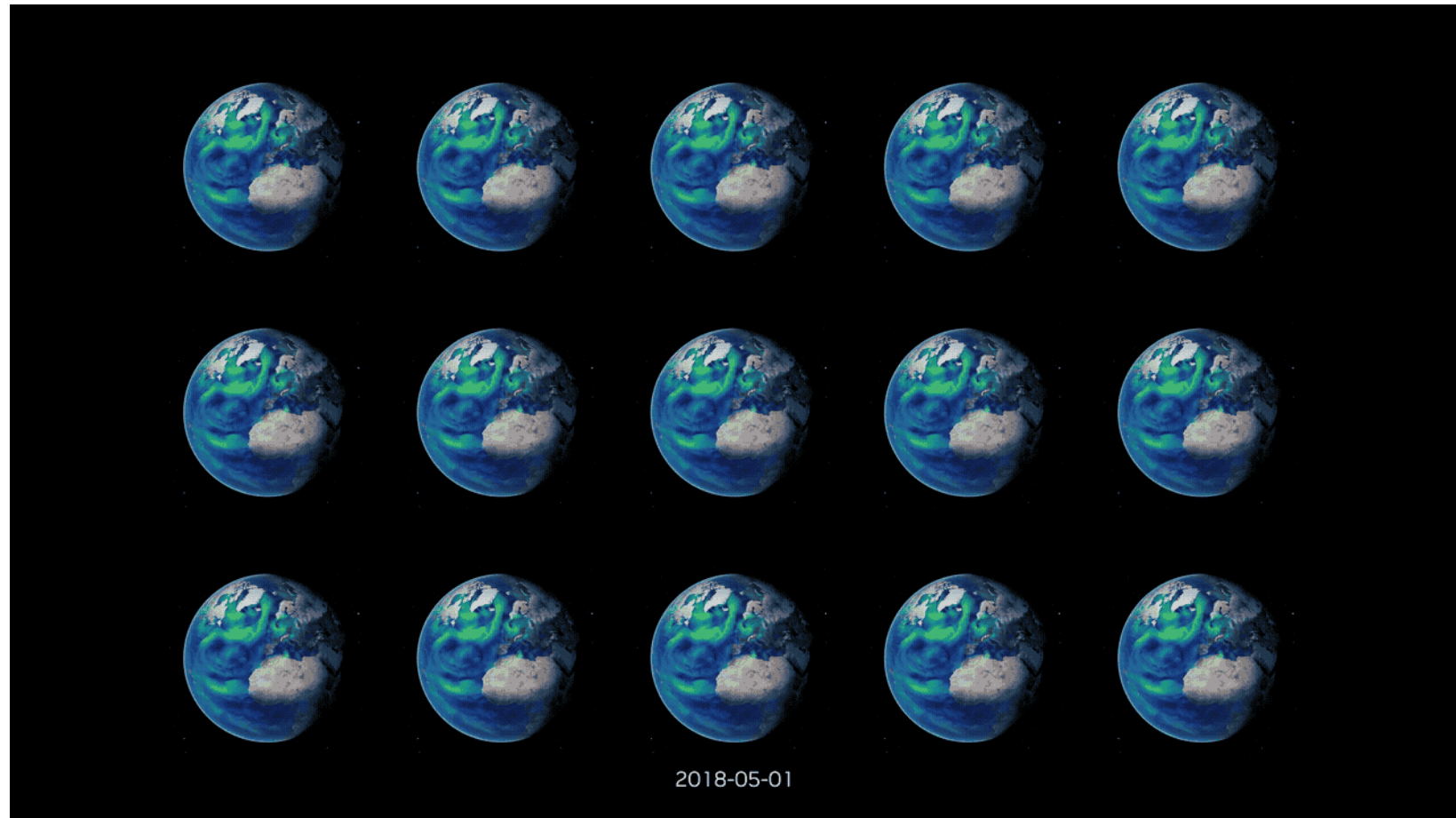


# Example: Forecasting

FourCastNet3



Ensemble



# Summary

- Data assimilation combines numerical models and observations to estimate the most accurate state of the Earth system (atmosphere, ocean, land, cryosphere).
- Forecast systems use methods like 3D-Var, 4D-Var, and Ensemble Kalman Filters to update model states based on observational data.
- These methods minimize the mismatch between model predictions and observations while accounting for uncertainties in both.
- AI is increasingly used in data assimilation to improve computational efficiency, integrate diverse data sources, and handle sparse or incomplete observations.