

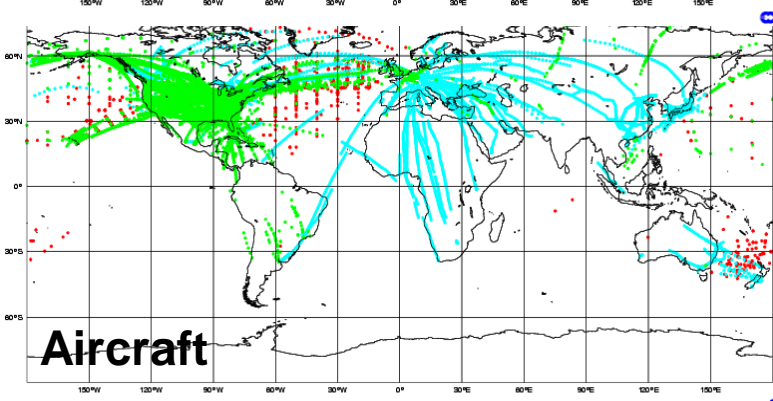
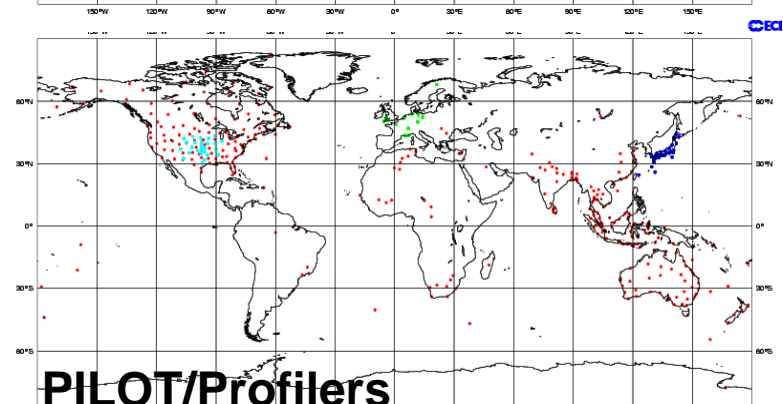
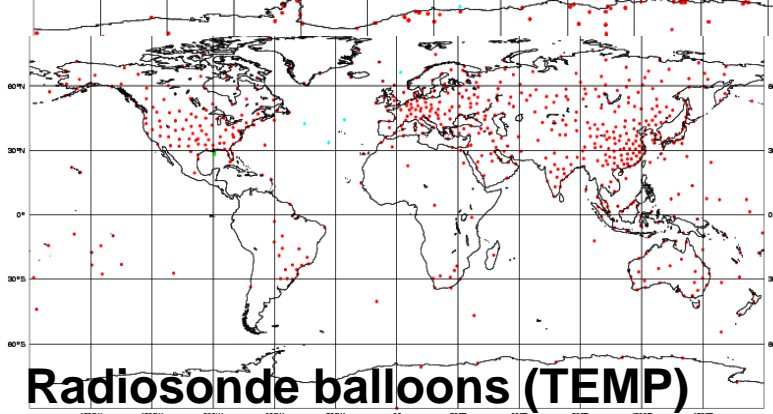
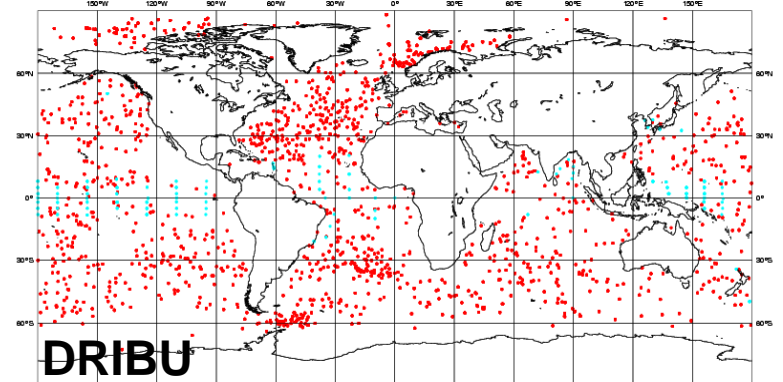
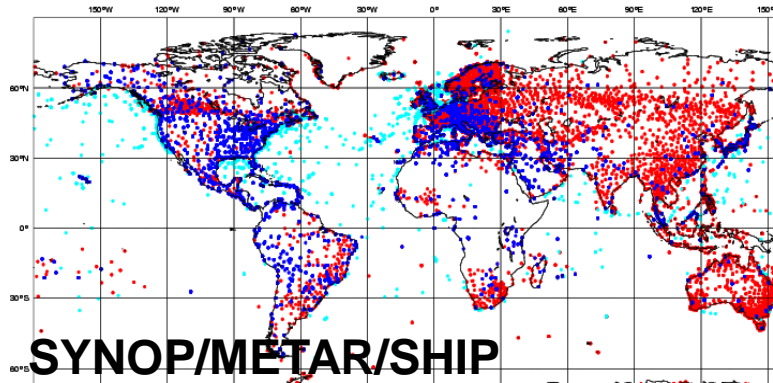
# Data Assimilation Training Course Summary and Q&A

# Data Assimilation

Data Assimilation has two main goals:

- Optimally blend information from **observations** and **model** to produce an accurate and physically consistent estimate of the **initial state** of the atmosphere and of the other components of the Earth System
- Quantify the **uncertainty** of our estimate of the initial state and sample this uncertainty in order to initialise an ensemble forecast

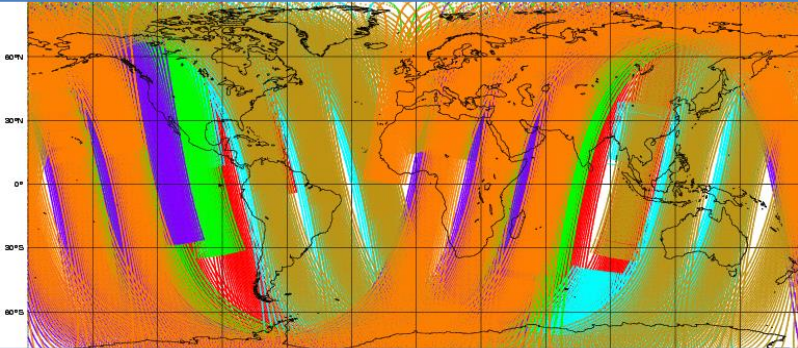
# Conventional/In situ observations



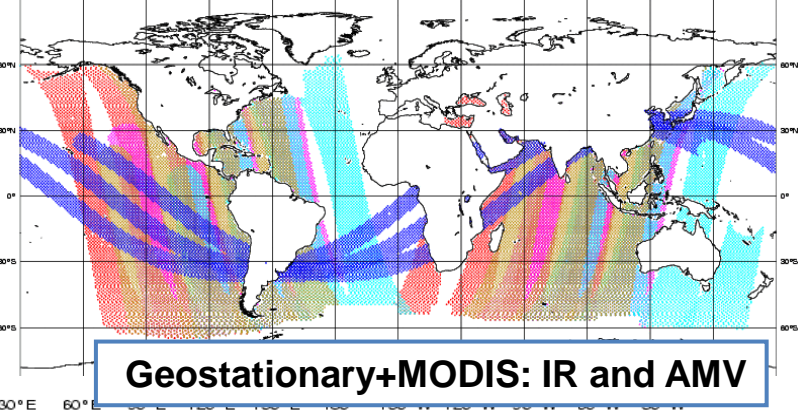
Sean Healy's talk

# Satellite observations

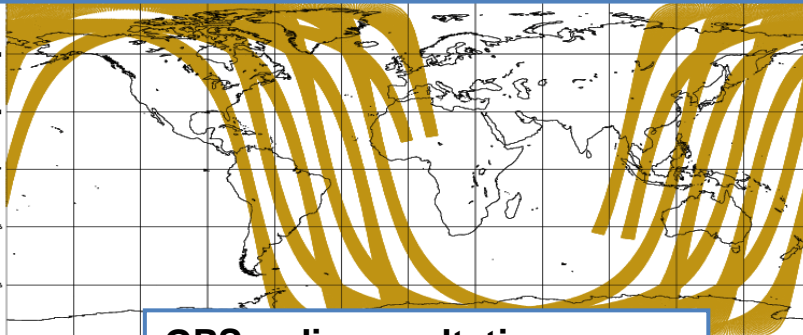
**Sounders: NOAA AMSU-A/B, HIRS, AIRS, IASI, MHS**



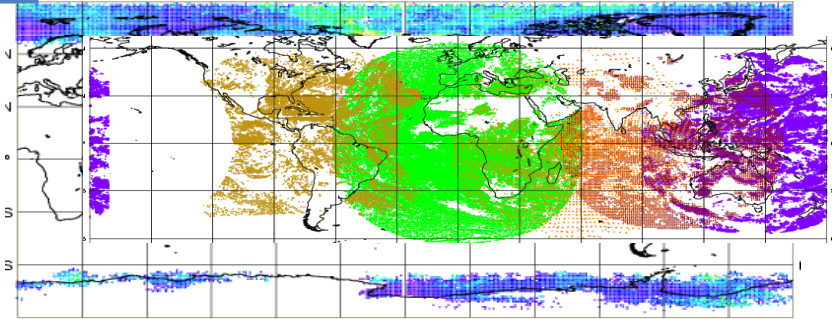
**Imagers: SSMI, SSMIS, AMSR-E,**



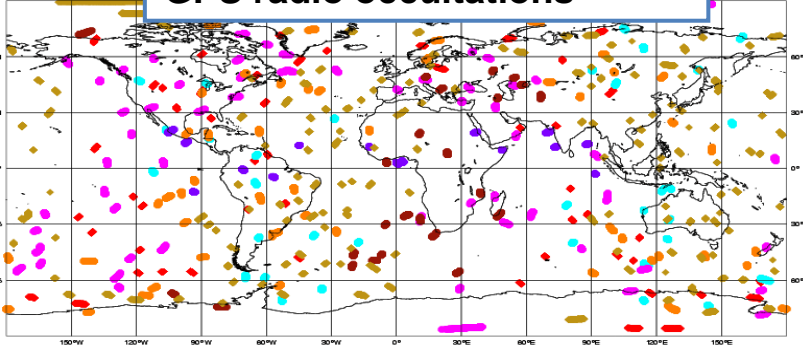
**Scatterometer ocean low-level winds: ASCAT**



**Geostationary+MODIS: IR and AMV**



**GPS radio occultations**



**Tony McNally and Sean Healy's talks**

# Observation errors

- Observations are affected by **errors** of different types

$$\mathbf{y} - \mathbf{y}^* = \varepsilon_o = \varepsilon_G + \varepsilon_M + \varepsilon_R + \varepsilon_H$$

- Observation errors are typically assumed to be zero-mean, Gaussian. When this assumption fails we need to do something special
- $\varepsilon_G$  (gross errors) are dealt with by **Observation Quality Control** techniques (Sean Healy's talk)
- Observations are assumed to be un-biased:

$$\langle \varepsilon_o \rangle = 0$$

- Biases are dealt with specific **Bias Correction** techniques: at ECMWF this is part of the analysis algorithm itself (e.g., **Variational Bias Correction**: Niels Bormann's talk)

# Observation errors

- In common DA algorithms we require not only the observations to be un-biased but also the background forecast to be un-biased:

$$\langle \varepsilon_b \rangle = 0$$

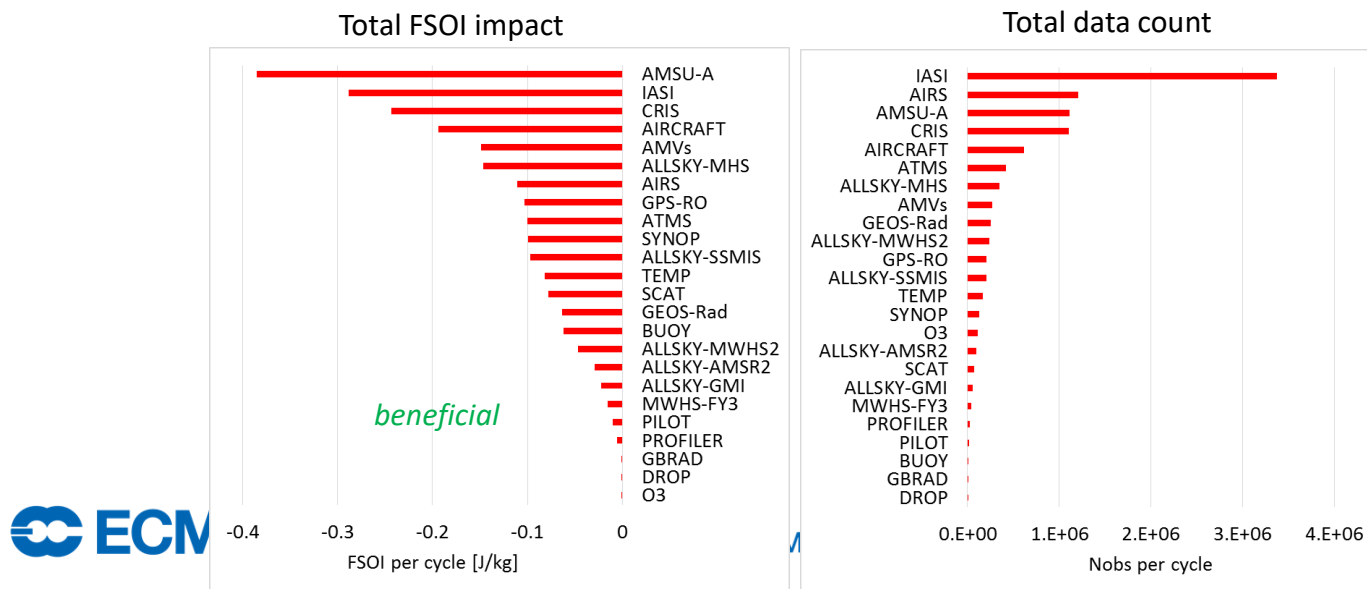
- But our only source of information about observation and forecast errors are observation departures:

$$y - H(x_b)$$

- We need to make further assumptions in order to disentangle observation and model error (Niels Bormann and Patrick Laloyaux's talks [on observation and model biases](#))

# Observation impact

- It is also important to monitor and evaluate the impact different types of observations have on the quality of the analyses and forecasts
- To do this we routinely look at [observation departures](#) (with respect to both analysis and forecast fields: see [Sean Healy's talk and practical sessions](#))
- We can also perform [Observing System Experiments \(OSEs\)](#) ([Alan Geer's talk](#))
- We routinely compute [adjoint-based diagnostic](#) of observation impact (Forecast Sensitivity to Observation Impact: [Alan Geer's talk](#))





# The forecast model

- A good model is able to effectively propagate information from past observations to the current analysis update => new batch of observations will only produce small corrections to the background => we are closer to the conditions of linearity of errors where current DA algorithms work best
- In [incremental 4D-Var](#) we not only require the full non-linear model to advance the state in time
- We also need its linearised versions ([Tangent Linear and Adjoint](#)) to propagate increments with respect to the first guess forward and backwards in time during the assimilation window (update of the linearised cost function and computation of its gradient with respect to the initial state)
- Developing and maintaining TL and ADJ codes is a complex task ([Philippe Lopez and Marcin Chrust's talks and practical](#)): but the availability of sophisticated TL and ADJ models is one of the main reasons for ECMWF success



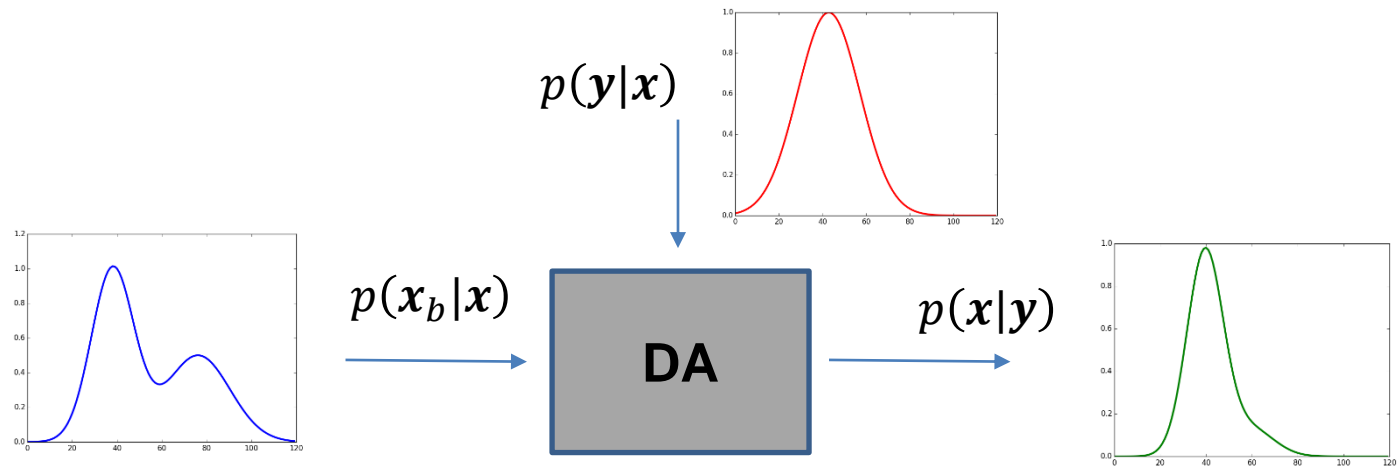
# Model errors

- Despite their increasing complexity and sophistication models are far from perfect!
- Many sources of model error: missing physical processes, errors in parametrizations of physical processes, discretisation errors (from continuous PDEs to discrete formulation), etc.,
- We represent model errors in two ways:
  1. **Stochastic errors**: explicitly perturbing the model integration in our ensemble data assimilation system (EDA; see Massimo Bonavita's talk – Assimilation Algorithms (5))
  2. **Model biases**: Using an explicit model error term in the 4D-Var cost function (weak constraint 4D-Var: see Sebastien Massart talk on 4D-Var and Patrick Laloyaux's talk on Model Error)

# Blending observations and model information: the Bayes perspective

# The Bayes perspective

- At an abstract level, we can think of the analysis process as updating our prior knowledge about the state, represented by a background forecast and the pdf of its errors, with new observations, represented by their values and the pdf of their errors:



$$p(x|y) = \frac{p(y|x)p(x)}{p(y)} = \frac{p(y|x)p(x_b|x)}{p(y)} \propto p(y|x)p(x_b|x)$$

- $p(x_b|x)$  = **prior pdf** (encapsulate our knowledge about the state before new observations)
- $p(y|x)$  = **observations likelihood** (pdf of the observations conditioned on the state)
- $p(x|y)$  = **posterior pdf** (updated pdf of the state after the analysis)
- $p(y)$  = **marginal pdf of the observations** (does not depend on  $x$ : normalising constant in Bayes' law)

# Particle Filters

$$p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x})p(\mathbf{x}_b|\mathbf{x}) \quad (1)$$

- In principle an analysis update requires being able to compute the product pdf of the random variables  $\mathbf{y}$ ,  $\mathbf{x}_b$ . This is usually not possible to do unless we choose very specific functional forms for the pdfs
- We thus need to make approximations
- One idea is to use Monte Carlo methods to sample and propagate the pdfs in (1) by an ensemble of states: [Particle Filters](#)
- This does not work for high dimensional systems as in NWP
- Need to make further assumptions on (1)

# Kalman Filter methods

- Need to make further assumptions on (1)
- Gaussian error pdfs => Gaussian posterior pdf

$$p(\mathbf{x}_a|\mathbf{y}) = \mathcal{N}(\mathbf{x}_a, \mathbf{P}^a)$$

$$\mathbf{x}_a = \mathbf{x}_b + \mathbf{K}(\mathbf{y} - \mathbf{H}(\mathbf{x}_b))$$

$$\mathbf{P}^a = (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{P}^b(\mathbf{I} - \mathbf{K}\mathbf{H})^T + \mathbf{K}\mathbf{R}\mathbf{K}^T = (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{P}^b$$

$$\mathbf{K} = \mathbf{P}^b\mathbf{H}^T(\mathbf{H}\mathbf{P}^b\mathbf{H}^T + \mathbf{R})^{-1} = \left( (\mathbf{P}^b)^{-1} + \mathbf{H}^T\mathbf{R}^{-1}\mathbf{H} \right)^{-1} \mathbf{H}^T\mathbf{R}^{-1}$$

- Solving directly these equations lead to Kalman Filter type DA methods: Optimum Interpolation, Kalman Filter, Extended KF, Ensemble KF (Massimo Bonavita's talk on KF and EnKF)
- These methods work well with low dimensional systems or small number of observations (O.I. in Snow analysis; Extended KF for soil moisture analysis, e.g. Patricia De Rosnay's talk on Land Data assimilation)
- For high-dim systems they require localisation: sophisticated localisation methods are needed to avoid losing information from non-local observations like satellite radiances

# Variational methods

- The Kalman Filter analysis update equation can be formulated as an equivalent minimization problem:

$$J(\mathbf{x}_0) = (\mathbf{x}_b - \mathbf{x}_o)^T (\mathbf{P}^b)^{-1} (\mathbf{x}_b - \mathbf{x}_o) + \sum_{t=0}^T (\mathbf{y}_t - H_t M_{0 \rightarrow t}(\mathbf{x}_0))^T \mathbf{R}_t^{-1} (\mathbf{y}_t - H_t M_{0 \rightarrow t}(\mathbf{x}_0))$$

- This is the basis of [Variational methods](#) (3D-Var, 3D-Var FGAT, 4D-Var: [see Sebastien Massart's lectures on 3-4DVar](#))
- Solving the KF update equation [globally](#) through iterative algorithms (conjugate gradient, Lanczos, Newton's methods)
- These methods do not require direct access to the elements of the error covariance matrices. We can represent [error covariances](#) by [operators](#) (i.e., pieces of code) acting on increments ([see Elias Holm talk on background error modelling](#))
- Variational methods work well on high dimensional systems and are generally used in global NWP

# Hybrid Data Assimilation methods

- The Kalman Filter equations require estimating and advancing in time not only the state but also its error covariance:

$$\mathbf{P}_t^a = (\mathbf{I} - \mathbf{K}\mathbf{H})\mathbf{P}_t^b (\mathbf{I} - \mathbf{K}\mathbf{H})^T + \mathbf{K}\mathbf{R}\mathbf{K}^T$$

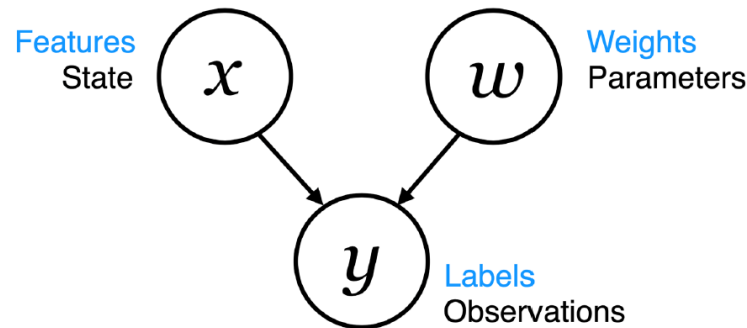
$$\mathbf{P}_{t+1}^b = \mathbf{M}\mathbf{P}_t^a \mathbf{M}^T + \mathbf{Q}_{t+1}$$

- 4D-Var can implicitly do this but only inside the assimilation window (12 hours at ECMWF)
- The idea of [Hybrid DA methods](#) is to combine a variational DA system to estimate the state with an ensemble data assimilation system (EnKF/EDA) to estimate and cycle the errors of the state (see Massimo Bonavita's talk on Hybrid data assimilation)
- Ensemble DA systems also provide the initial conditions for Ensemble Prediction
- All major global NWP Centres run Hybrid DA systems for Atmospheric DA



# Machine Learning and Data Assimilation

- We have also learned this week how Data Assimilation can be extended to estimate not only the initial state but also the model which underlies the evolution of the state (Alan Geer's talk)



- **Data Assimilation provides a broader theoretical framework than standard Machine Learning!** However we can take advantage of the many Machine Learning software tools and concepts and use them in our DA applications!

# Earth System Data Assimilation

- We have discussed Data Assimilation methods with an emphasis on global NWP
- The DA methods presented are however general: which one to apply to a given problem depends on the characteristics of the problem (size of the state vector, number and quality of observations, available computing resources, available manpower,...)
- You have seen applications in [Atmospheric Composition DA](#) (4D-Var: [Richard Engelen's talk](#)); in [Ocean Data Assimilation](#) (3D-Var FGAT: [Hao Zuo's talk](#)); in [Land Data Assimilation](#) (O.I., Simplified EKF: [Patricia de Rosnay's talk](#))
- In current ECMWF DA the Earth system's components are only weakly coupled (through a coupled model background forecast)
- [Phil Browne's talk](#) has given you our current perspective of some of the challenges and the potential benefits of a stronger [coupling in the data assimilation](#) for the different components of the Earth System

# Earth System Data Assimilation

- We have discussed Data Assimilation methods for the Earth System with an emphasis on producing the best initial state estimate for [forecasting at short, extended and seasonal timescales](#)
- An increasingly important application of Earth System DA is to help to [reconstruct the past climate and weather](#) (see Dinand Shepers' talk on Reanalysis methods)
- As DA methods have dramatically improved over the years we are able to make better use of past observational records and more robustly estimate climatic trends

# Earth System Data Assimilation

- Our main goal this week has been to provide you with a structured introduction to DA methods for Earth System and an overview of the main application areas
- Where to look for material about current topics and challenges:
  - ECMWF Annual Seminar 2018: Earth System Assimilation (<https://www.ecmwf.int/en/learning/workshops/annual-seminar-2018>)
  - ECMWF Annual Seminar on Observations 2021 (<https://events.ecmwf.int/event/217/>)
  - ECMWF Annual Seminar on Reanalysis 4-8 September 2023
  - ECMWF-ESA Workshops on Machine Learning for Earth System Observation and Prediction (2020/21/22; search online for presentations), next edition at ESA-ESRIN (Rome), May 2024
  - International Symposium on Data Assimilation (ISDA): yearly in person event next edition in Bologna (Italy, 16-20 Oct 2023; <https://eventi.unibo.it/isda2023>)
  - International Symposium on Data Assimilation (ISDA, online, running event: <https://isda-online.univie.ac.at/>)

Thank you for being such a  
receptive and engaged audience!

Questions?