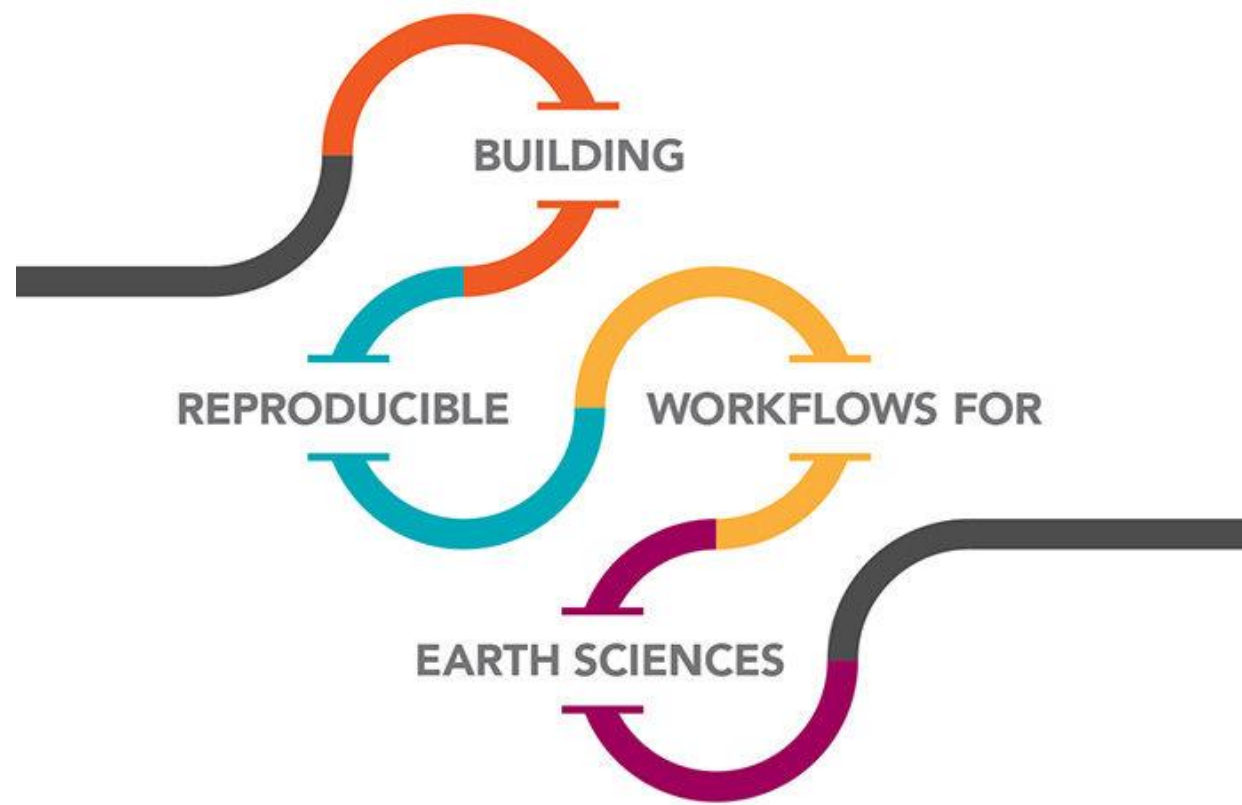# Reproducible workflows – setting the scene
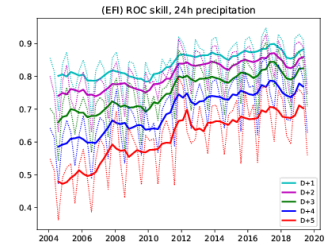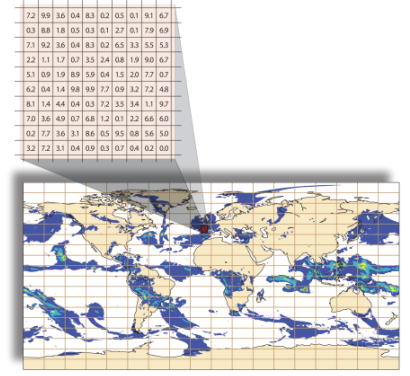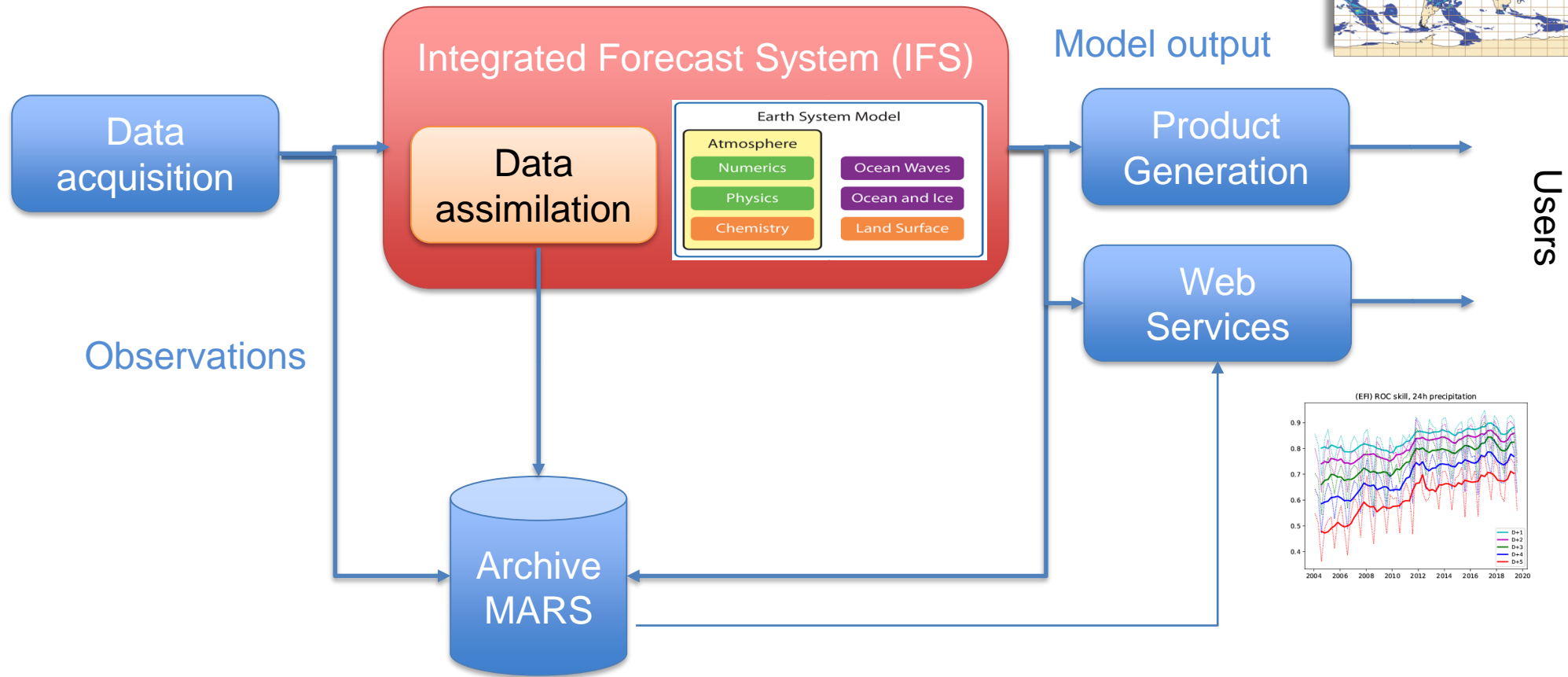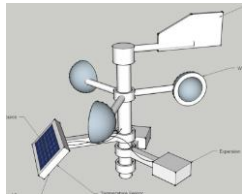
## Why are we here?

Stephan Siemen

Development Section, Forecast Department

# What does ECMWF do?



Integrated Forecast System (IFS)

Data acquisition

Data assimilation

Earth System Model

Atmosphere
- Numerics
- Physics
- Chemistry

- Ocean Waves
- Ocean and Ice
- Land Surface

Model output

Product Generation

Web Services

Users

Observations

Archive MARS
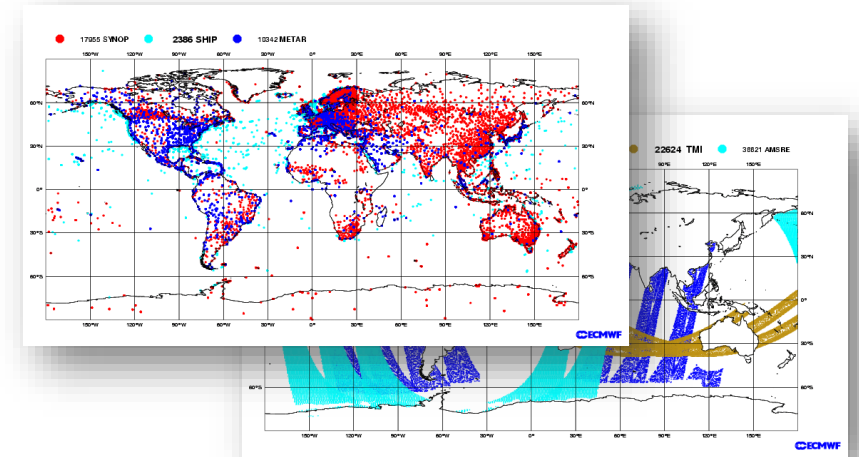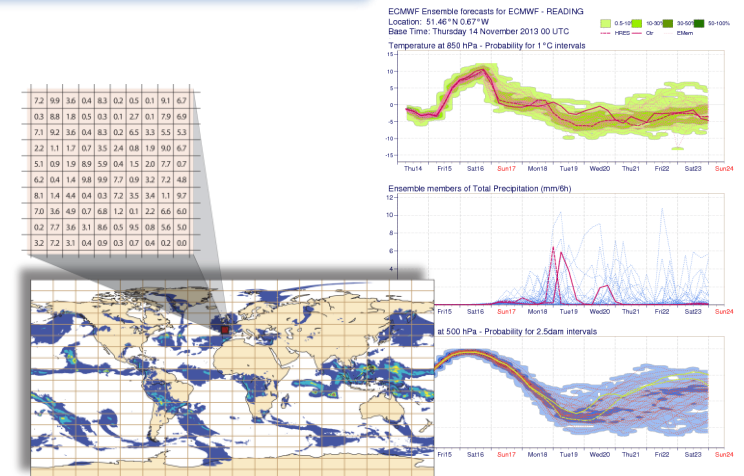
# The three stages of a research workflow



1. data acquisition & cleaning & filtering

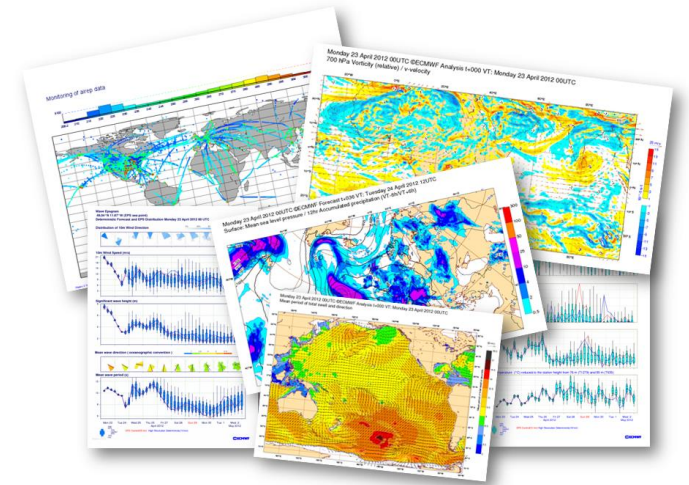2. data processing – running models

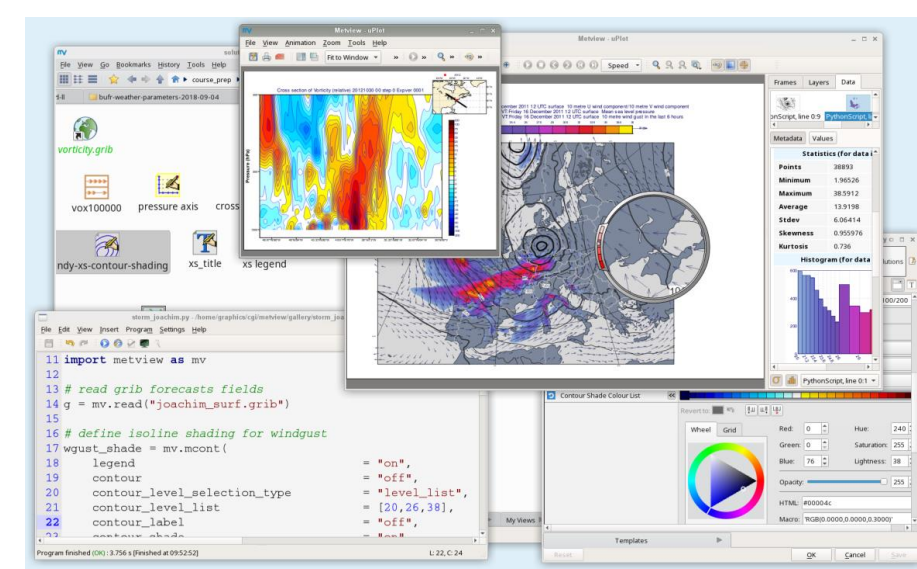3. data analysis – presenting outcomes

# What is reproducibility anyway?

- Not only getting the same result
  - Hard to achieve anyway with different architectures and environments

- Having a recipe / procedure to follow to get to the same result
  - Allow to evaluate changes to input data and algorithm on output → do science

- Having workflow documented allows sharing and scrutiny
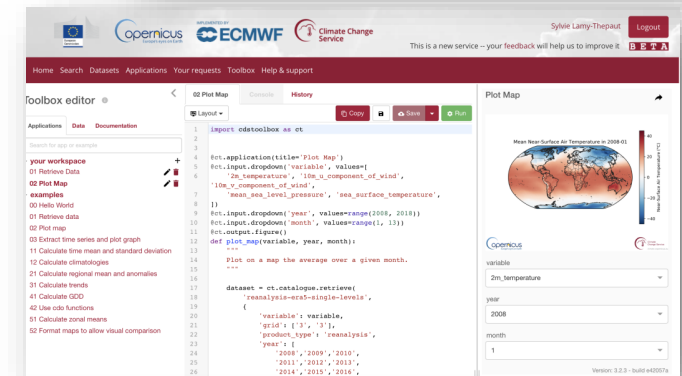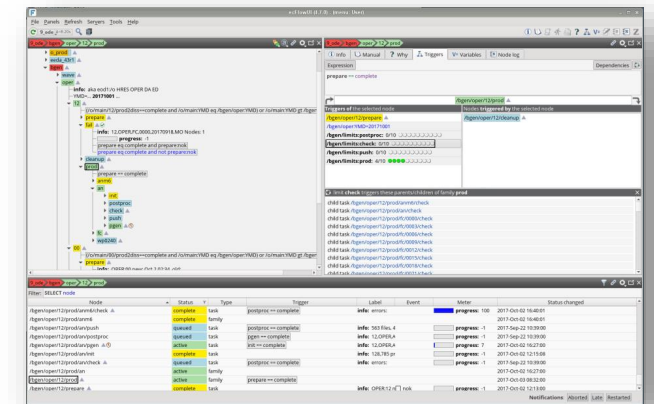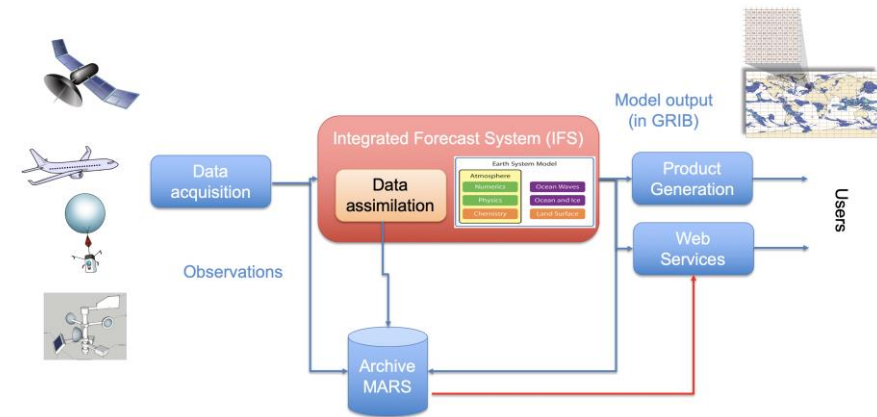
# Why is it important?



- Important proof that science is solid
  - No one believes a result if it constantly changes
  - With changing input a hypothesis can be confirmed/corrected

- Allows easy sharing of scientific work and allows others to follow and reflect on work

- Big Data workflows need to carefully designed because of their high costs
  - Data amounts are challenging to handle and prone to change
  - Manual work does not scale and needs automation
  - This is especially true for AI / Machine Learning applications

# Why is it important to ECMWF?

- Reproducibility has always been important to drive innovation on model developments

- As an operational centre users rely on ECMWF to provide robust results in a controlled and stable environment

- With data amounts continuing to grow complex processing jobs need to be moved to the data
  - Require flexible and high-level workflows

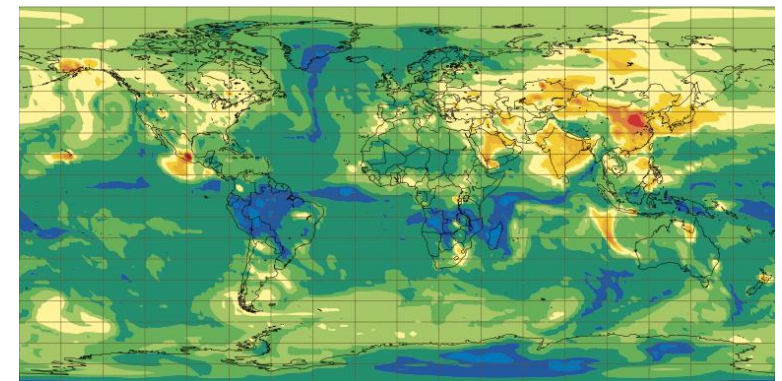- Being a central part of a large scientific community, ECMWF is keen to share workflows and offer training

# EU Copernicus programme



- Large programme to make earth observations and monitoring data (freely) available

- Data is validated and from verified sources

- Sponsored by the European Commission

- ECMWF is operating two services on behalf of the EC


- Copernicus Atmospheric Monitoring Service (CAMS)

  – Air pollution, chemicals, sand & dust

  – https://atmosphere.copernicus.eu


- Copernicus Climate Change Services (C3S)

  – Reanalysis data, seasonal forecasts, climate scenarios

  – https://climate.copernicus.eu



**ECMWF** EUROPEAN CENTRE FOR MEDIUM-RANGE WEATHER FORECASTS

*Total column of sulphur dioxide*
*(provided by the Copernicus Atmosphere Monitoring Service)*

# So what have we done?

# What makes a workflow reproducible?

# General advice

- Document all data, files, and operations that occur on data and files. Keep meta data of all components and operations

- Create the workflow as a sequence of small steps that are connected. Use intermediate outputs from a step input for the next step

- Automate all work as much as possible, and avoiding manual intervention in the workflow – let the automation (scripts) be the documentation

- Key are clear interfaces between components - APIs

# Technical developments of recent years help

- Git made version control main stream
  - Big ecosystem of tools → GitHub enables community for developers

- Python and its large eco system of packages

- Conda & pip help to define environments
  - Let's user easily recreate environments

- Jupyter notebooks made it popular to document codes
  - Easy to share interactive environments

- Containers make it easier to "freeze" environments
  - Next steps up are Kubernetes and service meshes to define whole services

# We have an interesting workshop in front of us

- Three days of great talks and demonstrations

- Discussions to exchange experiences