

# **Clustering Techniques and their applications at ECMWF**

Laura Ferranti

European Centre for Medium-Range Weather Forecasts

# Outline

Cluster analysis – introduction

Clustering products at ECMWF

Flow dependent verification

Predictability of Euro-Atlantic regimes from medium to seasonal ranges

A 2-dim space framework to detect early warnings of cold spell over Europe

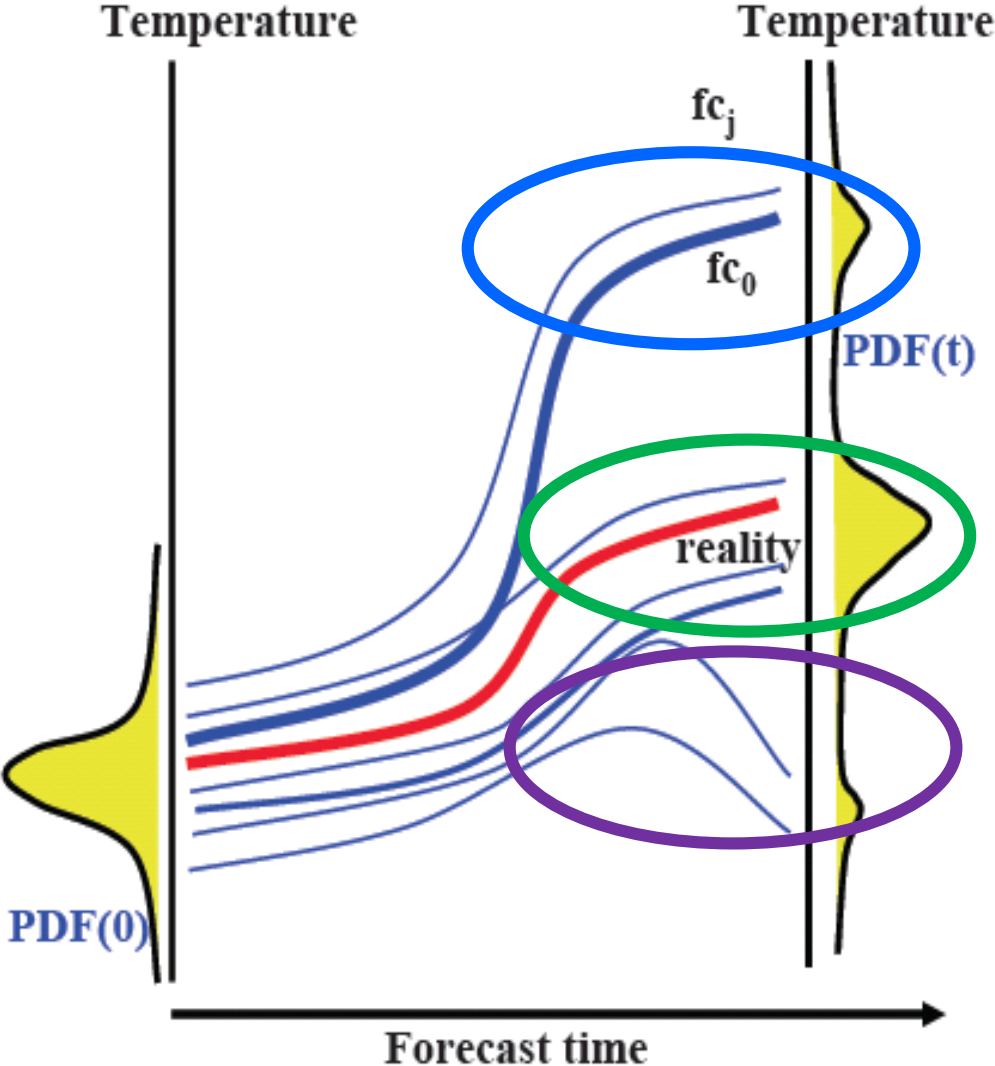
## Cluster analysis - Generalities

“Cluster analysis deals with **separating data into groups whose identities are not known in advance**. In general, even the “correct number” of groups into which the data should be sorted is not known in advance.” *Daniel S. Wilks*

### **Examples of use of cluster analysis in weather and climate literature:**

- Grouping daily weather observations into synoptic types (Kalkstein et al. 1987)
- **Defining weather regimes from upper air flow patterns** (Mo and Ghil 1998; Molteni et al. 1990)
- **Grouping members of forecast ensembles** (Tracton and Kalnay 1993; Molteni et al 1996; Legg et al 2002)

# Example – Grouping members of Forecast Ensembles



## Clustering techniques:

- Exclusive Clustering - data are grouped in an exclusive way
- Overlapping Clustering - fuzzy set of clusters data

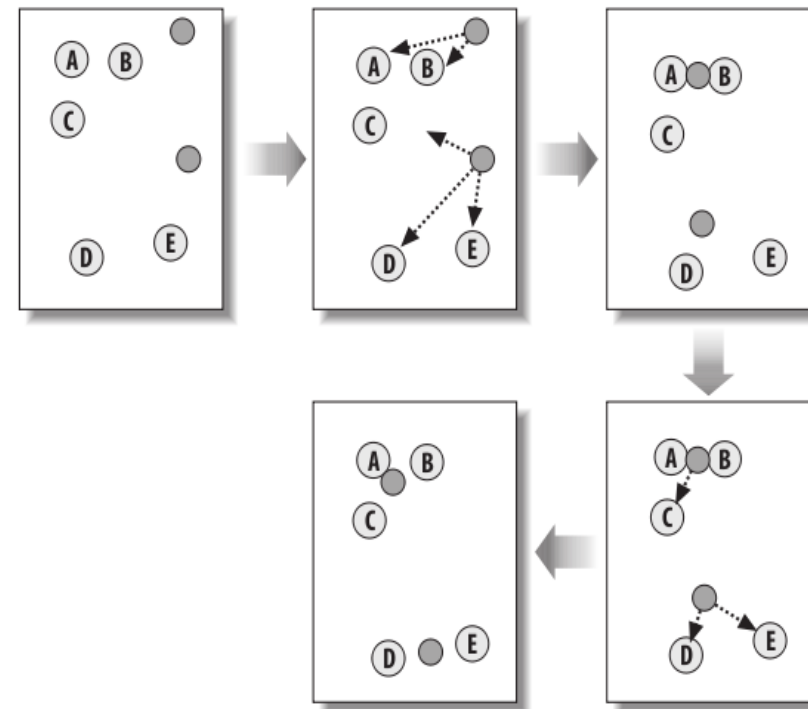
A widely used exclusive clustering approach is called K-means method. K is the number of clusters into which the data will be grouped ([this number must be specified in advance](#)).

## Cluster analysis - K-means method

➤ For a given number  $k$  of clusters, the optimum partition of data into  $k$  clusters is found by an algorithm that takes an initial cluster assignment (based on the distance from random seed points), and iteratively changes it by assigning each element to the cluster with the closest centroid, until a “stable” classification is achieved. (A cluster centroid is defined by the average of the PC coordinates of all states that lie in that cluster.)

➤ This process is repeated many times (using different seeds), and for each partition the ratio  $r_k^*$  of variance among cluster centroids (weighted by the population) to the average intra-cluster variance is recorded.

➤ The partition that maximises this ratio is the optimal one.



## Cluster analysis - How many clusters?

The need of specifying the number of clusters can be a disadvantage of K-means method if we don't know in advance what is the best cluster partition of the data set in question. However there are some criteria that can be used to choose the optimal number of clusters.

- **Significance:** partition with the highest significance with respect to predefined Multinormal distributions
- **Reproducibility:** We can use as a measure of reproducibility the **ratio of the mean-squared error of best matching cluster centroids from a N pairs of randomly chosen half-length datasets** from the full actual one. The partition with the highest reproducibility will be chosen.
- **Consistency:** The consistency can be calculated both with respect to variable (for example comparing clusters obtained from dynamically linked variables) and with respect to domain (test of sensitivities with respect to the lateral or vertical domain).

# Clustering product at ECMWF:

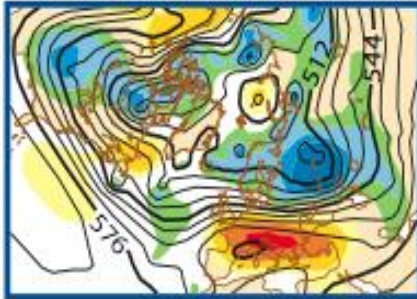
1. Identification of cluster scenarios to reduce the dimension of the ensemble forecast distribution (51 → max of 6)
2. Association of each cluster scenario to climatological weather regime.



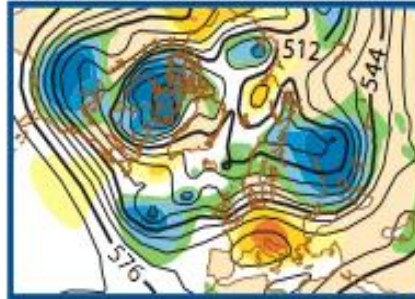
# Regime transitions within a time window

Day 5 to day 7 - 9 February 2011 – 3 scenarios 2 possible transitions

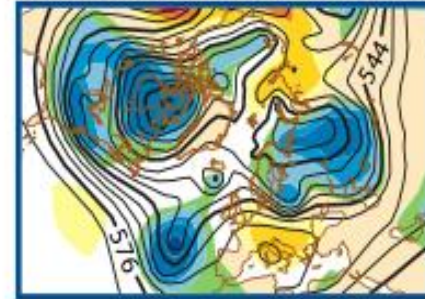
Population: 22. Representative member: 0



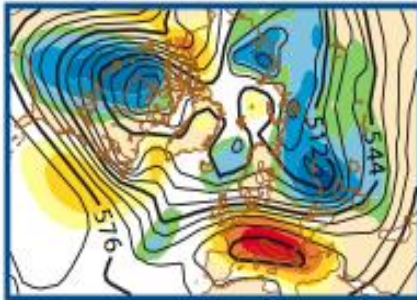
Population: 22. Representative member: 0



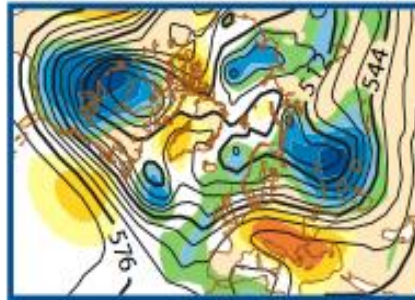
Population: 22. Representative member: 0



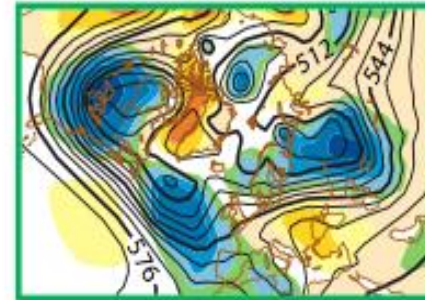
Population: 15. Representative member: 29



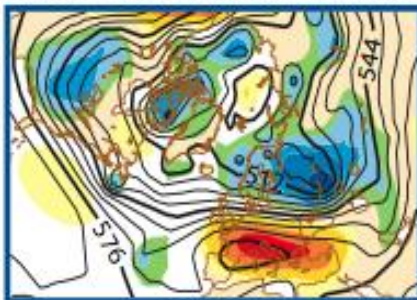
Population: 15. Representative member: 29



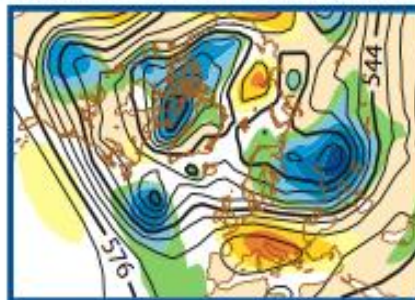
Population: 15. Representative member: 29



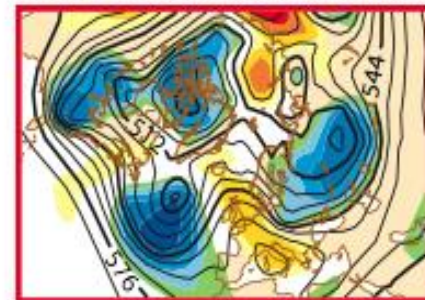
Population: 14. Representative member: 46



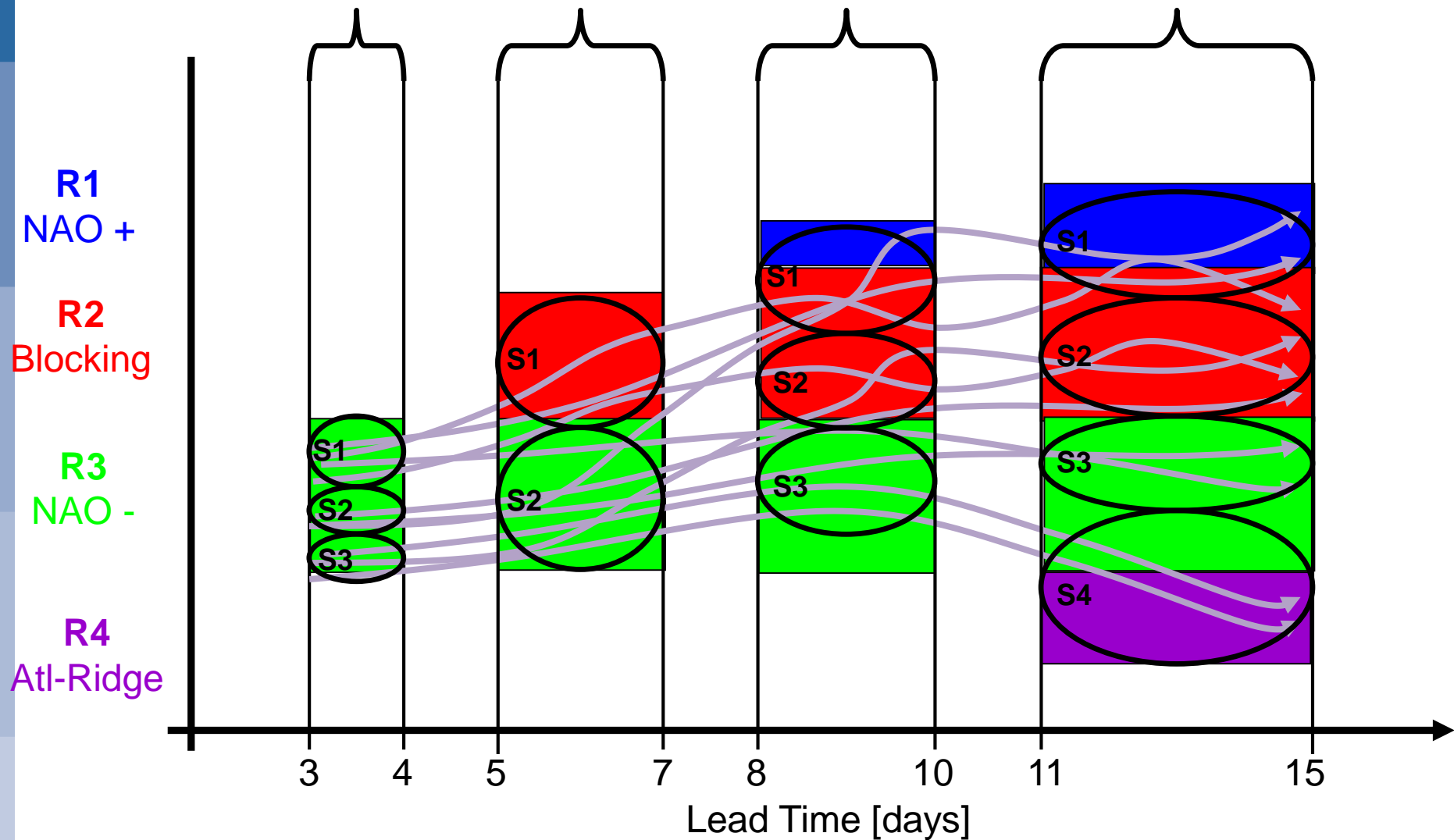
Population: 14. Representative member: 46



Population: 14. Representative member: 46



# Regimes & Scenarios



# Cluster product at ECMWF: large scale climatological regimes

**To put the daily clustering in the context of the large-scale flow** and to allow the investigation of regime changes, the new ECMWF clustering contains **a second component**. Each cluster is attributed to one of a set of four pre-defined climatological regimes.

**NAO+:** Positive phase of the North Atlantic Oscillation.

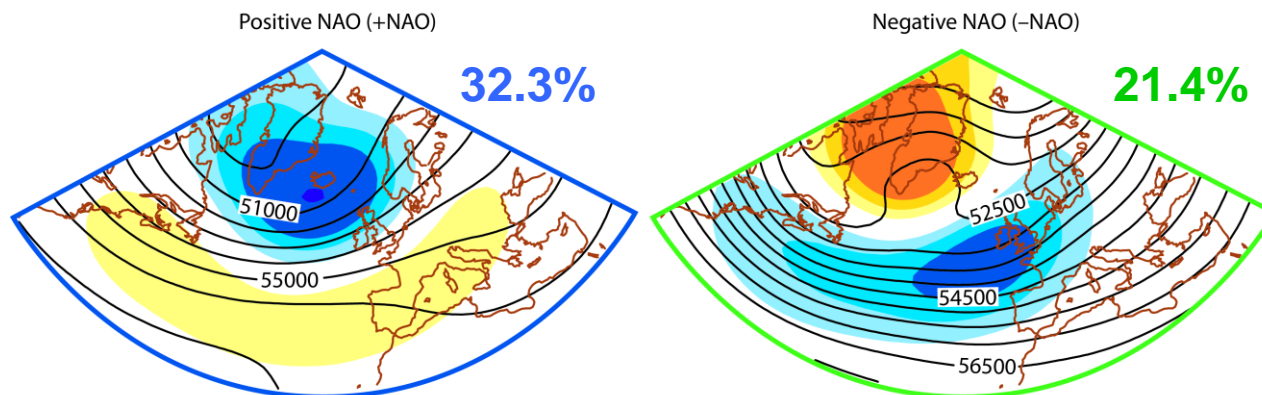
**BLO:** Euro-Atlantic blocking.

**NAO- :** Negative phase of the North Atlantic Oscillation.

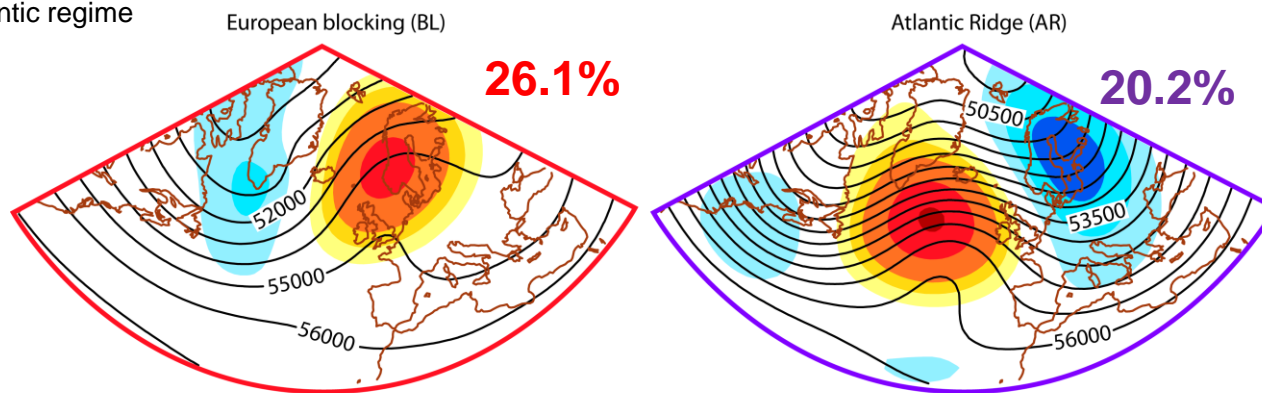
**AR:** Atlantic ridge.

# Regimes based on clustering of daily anomalies for 29 cold seasons ( October to March 1980-2008)

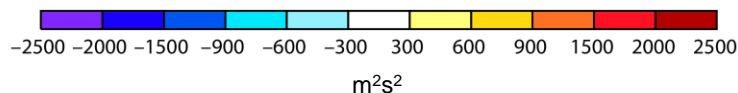
500 hPa geopotential



- Obtain well-known Euro-Atlantic regime patterns



'k means' clustering applied to EOF pre-filtered data (retaining 80% of variance)



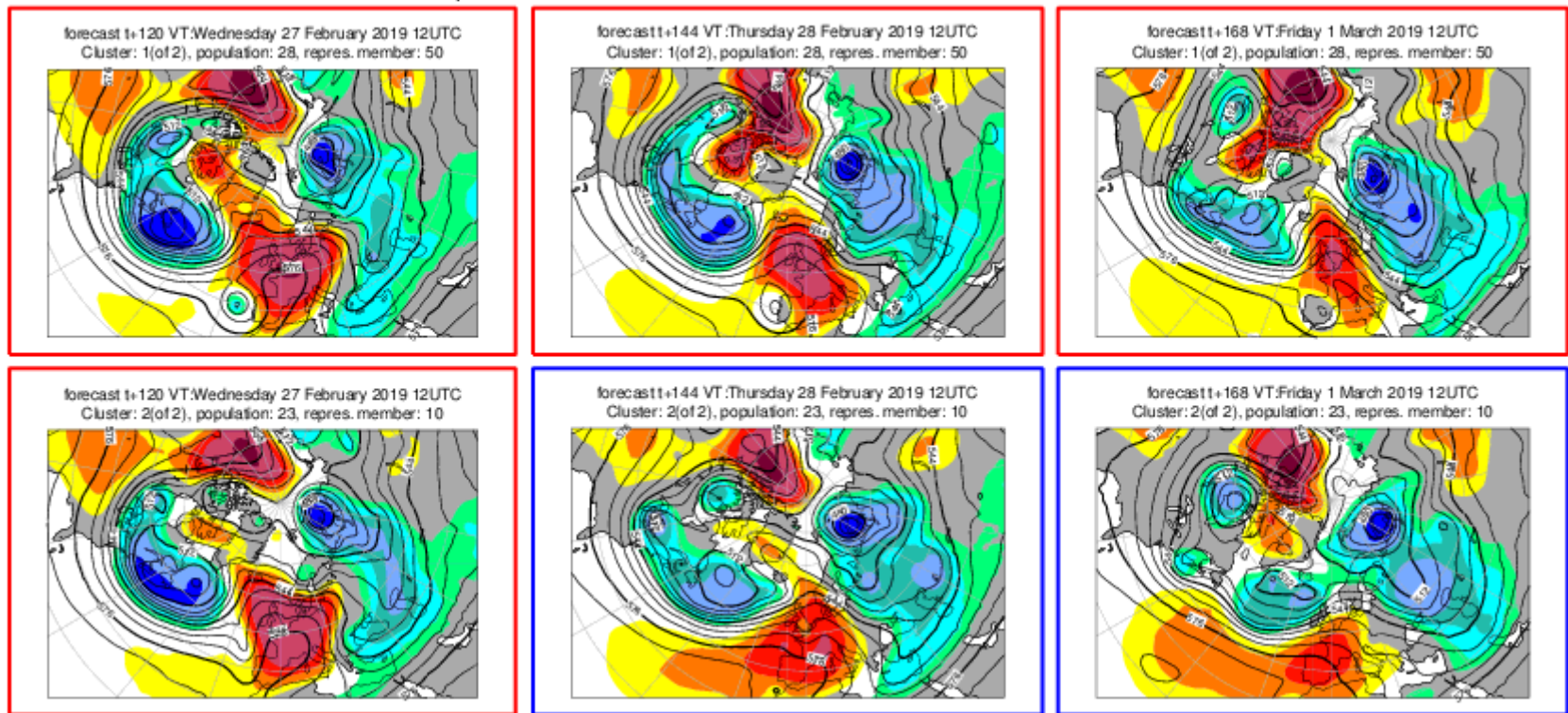


# Cluster scenario



## 2 Cluster scenarios 1 possible transition

Friday 22 February 2019 12UTC ECMWF EPS Cluster scenario - 500 hPa Geopotential  
 Reference step t+120-168 Domain 75/340/30/40 Cont. in cluster=1 Det. in cluster=2



+5days

+6 days

+7 days

# Cluster product at ECMWF:

The ECMWF clustering is one of a range of products that summarise the large amount of information in the Ensemble Prediction System (EPS).

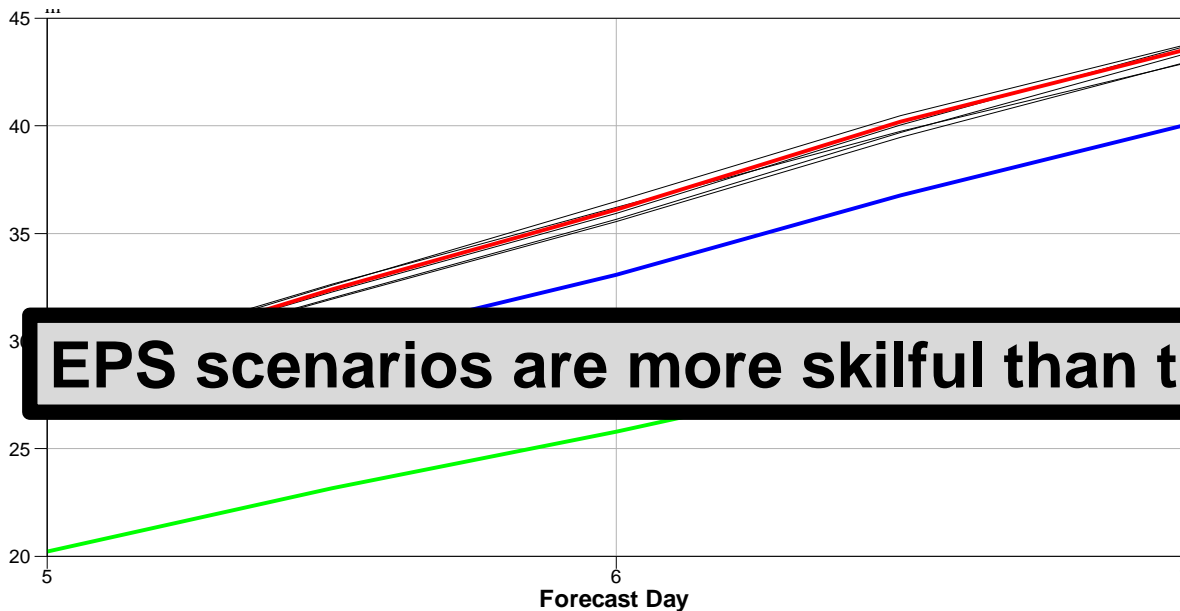
The clustering gives an overview of the different synoptic flow patterns in the EPS. The members are grouped together based on the similarity between their 500 hPa geopotential fields over the North Atlantic and Europe.

They are archived in MARS and available to forecast users through the operational dissemination of products.

A graphical clustering product is available for registered users on the ECMWF web site: <http://www.ecmwf.int/en/forecasts/charts/>

# What is the performance of the most probable scenarios?

Scenario distribution —  
Full EPS (50 members) —  
Reduced EPS —  
Ensemble Mean —



A randomly chosen 6 member ensemble has a CRPS equivalent to that of the ensemble mean.

Probabilistic scores depend largely on the ensemble size.

A large ensemble provides a more detailed and more reliable estimate of the forecast distribution.

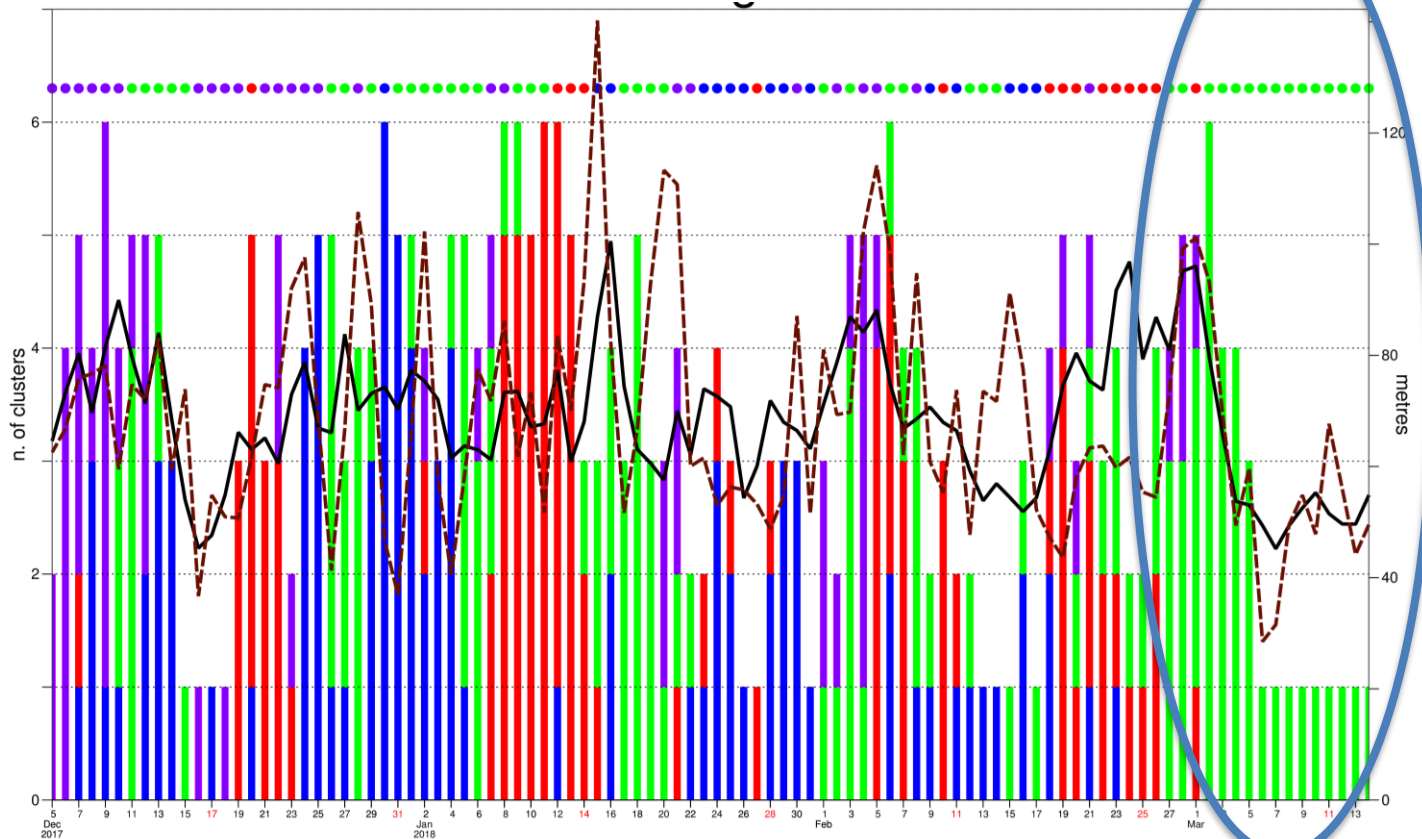
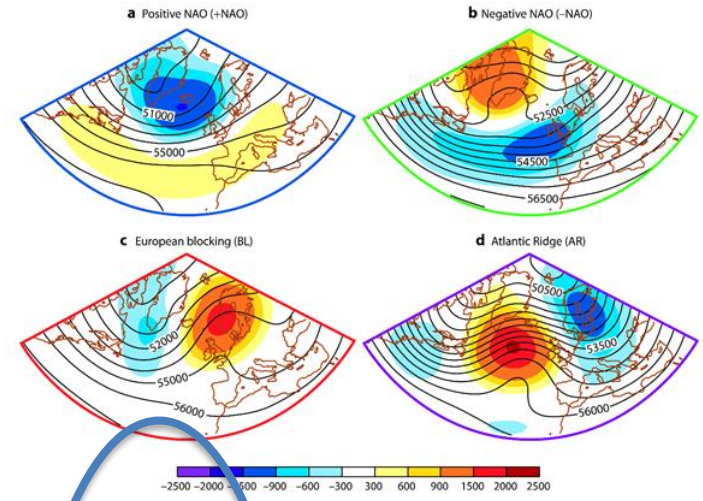
**EPS scenarios are more skilful than the ensemble mean**

# Assessing cluster scenarios performance :

## December 2017 – March 2018

Forecast range 7days

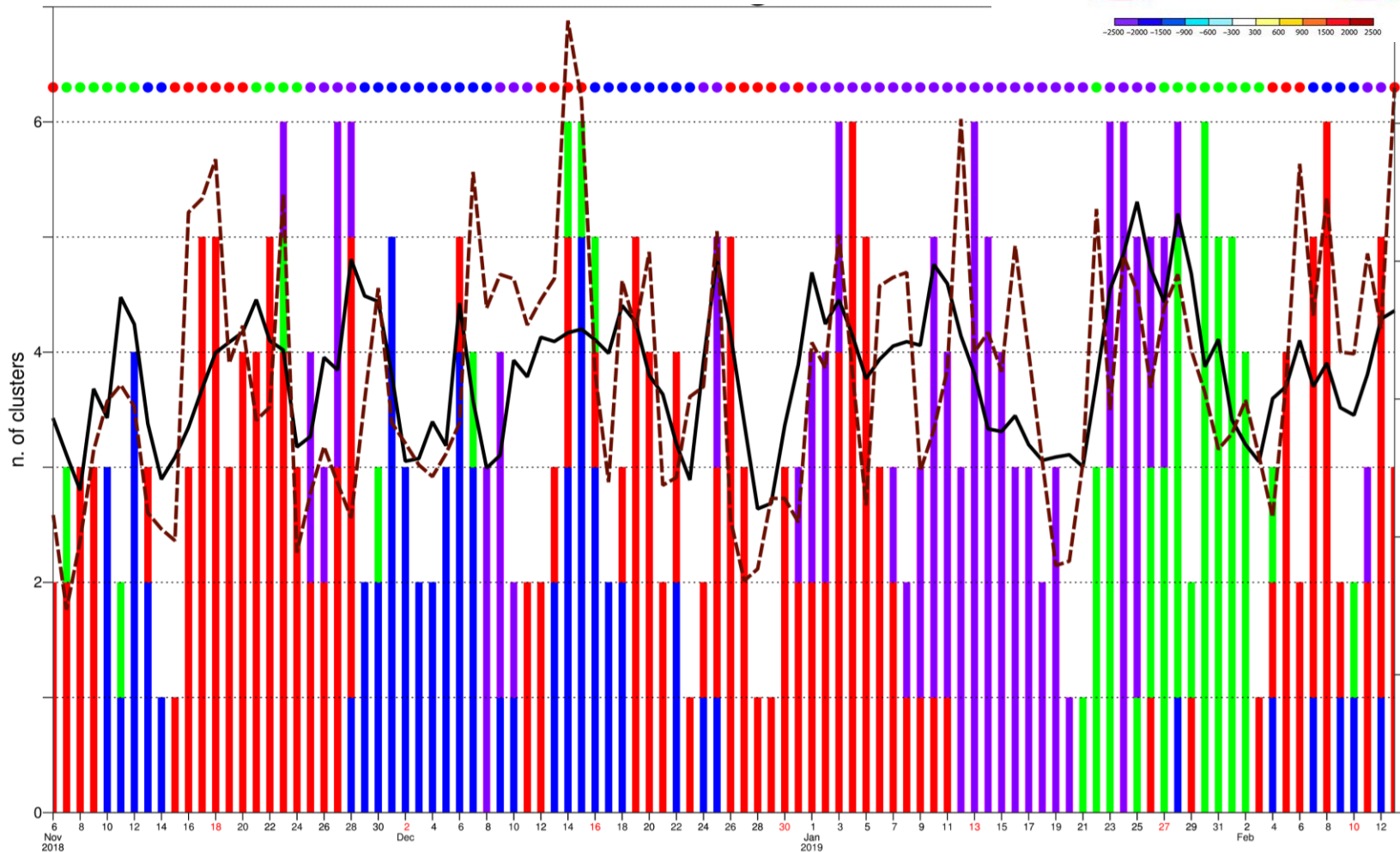
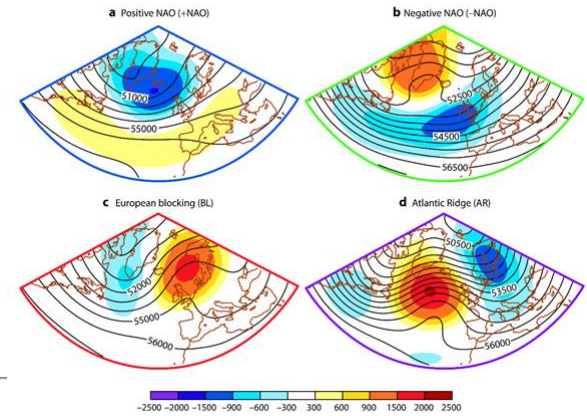
Spread error  
 —————  
 - - - - -





# Assessing cluster scenario performance:

November – Feb 2019:  
forecast range Day 7



The cluster product provides the users with a set of weather scenarios that appropriately represent the ensemble distribution

The classification of each EPS scenario in terms of pre-defined climatological regimes provides an objective measure of the differences between scenarios in terms of large-scale flow patterns. **This attribution enables flow-dependent verification** and a more systematic analysis of EPS performance in predicting regimes transitions

This clustering tool can be used to create EPS clusters tailored to the users' needs (e.g. different domain, different variables)

# Flow dependent verification over the Atlantic sector

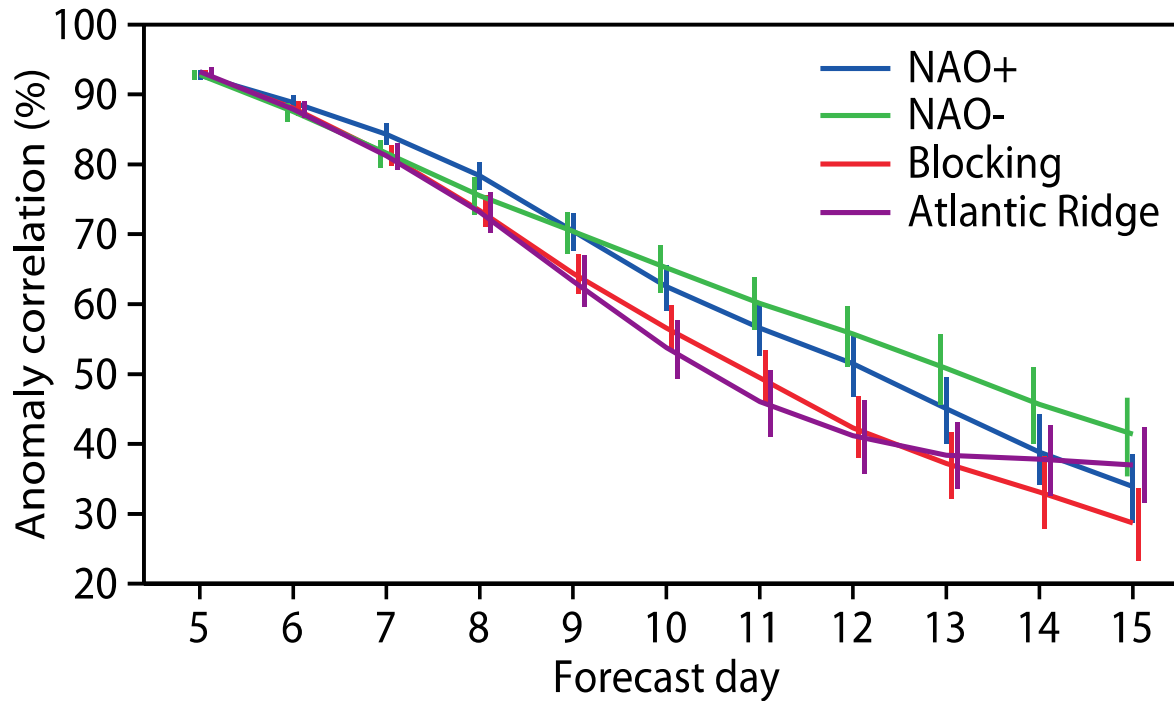
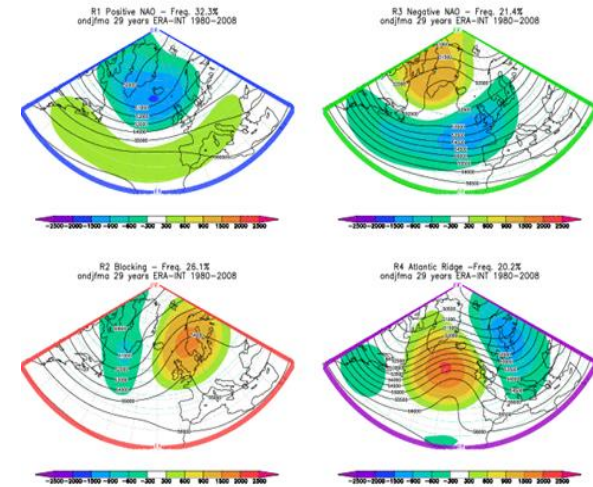
Identifying the flow configurations that lead to a more/less accurate forecast and quantifying the skill changes.

The concept of weather regimes is used to classify different flow configurations.

Oper. Forecast data: ENS cold season (Oct to April)  
2007-2012 operational analysis

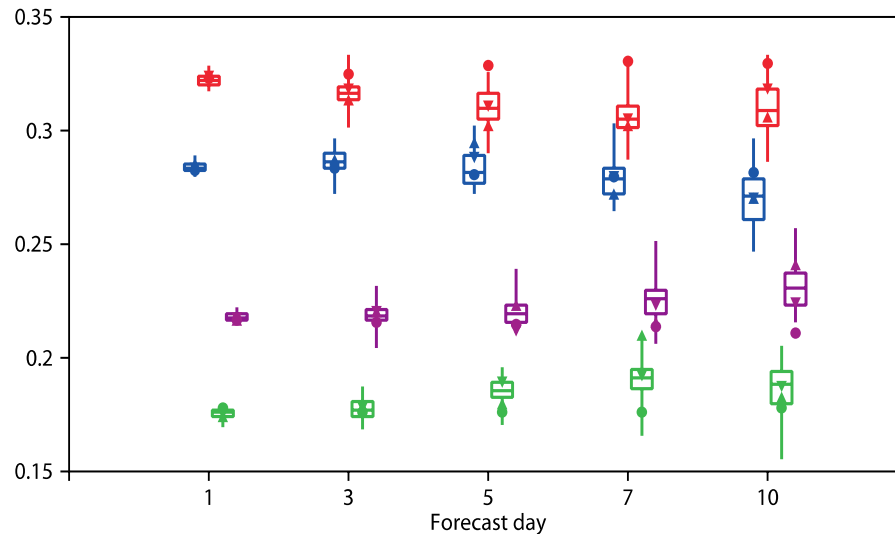
*( Ferranti et al. 2015 QJRMS)*

# Which flow pattern leads to a more/less accurate forecasts?



Anomaly correlation of the ensemble means for the four forecast categories as a function of forecast range. The bars, based on 1000 subsamples generated with the bootstrap method, indicate the 95% confidence intervals.

# Climatological frequency distribution for the 4 Euro-Atlantic regimes as simulated by the ECMWF ensemble at different forecast ranges



Red indicate the frequency of the BL regime, blue (green) the frequency of the NAO+ (NAO-) and violet the frequency of the AR regime. The observed frequencies are indicated by a circle while the frequencies from the ECMWF operational high resolution and the unperturbed forecasts are indicated by a pointing down and a pointing up triangle respectively.

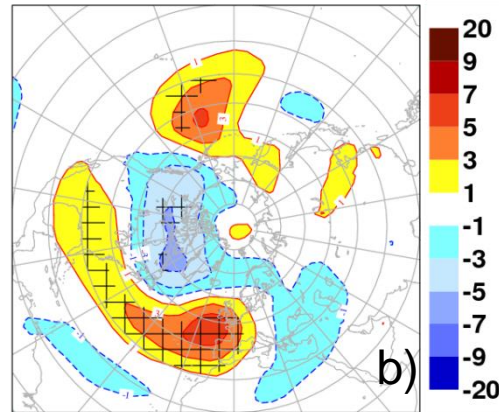
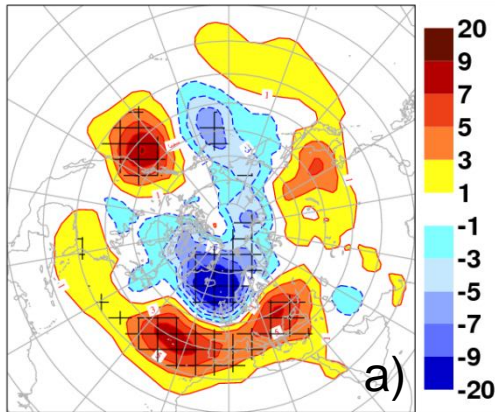
# Poor forecasts at day 10

The performance of the Ensemble is assessed by stratifying the cases according to their initial conditions as well as their accuracy at forecast day 10.

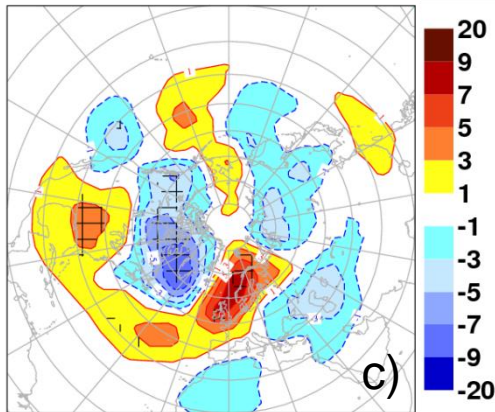
Poor (good) forecasts => RMSE of the ensemble mean larger (smaller) than the upper (lower) fifth of the whole RMSE distribution.

The RMSE is computed over the European domain at day 10. For each group and each category we compute composites maps of z500 anomalies at several time steps.

# Forecasting regimes transitions:



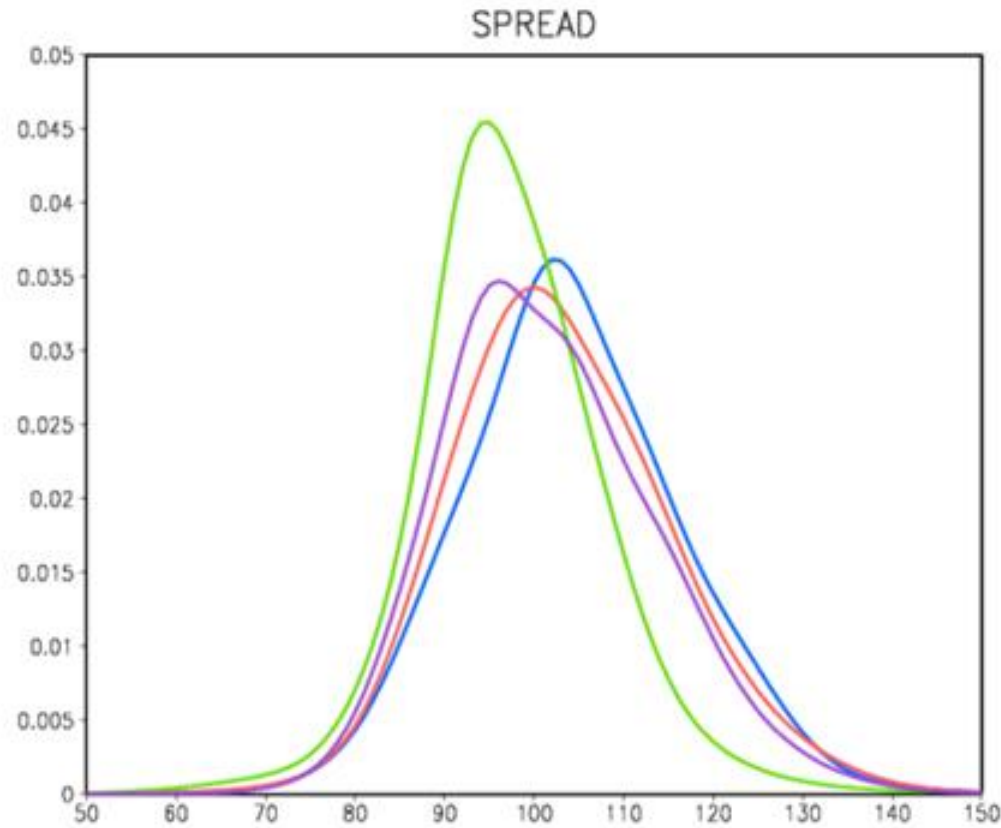
Composites of z 500 anomalies for all the forecasts initiated with flow configuration close to the NAO+ and with a RMSE at day 10 exceeding the upper quintile of the RMSE distribution.



- a) anomaly composites at the initial conditions;
- b) anomaly composites for the forecasts at day 10;
- c) anomaly composites of the corresponding verifying analysis. Hatched shading indicates statistical significance at the 10% level

	Day 0	Day 1	Day 5	Day 7	Day 10
<b>Forecasts with large RMSE at day 10</b>					
NAO+	100	81	56, 44	54, 40	37, 21
BL	0	8	28, 40	35, 53	42, 51
NAO-	0	2	0	2	2, 5
AR	0	9	16	9, 5	19, 23

NAO+ (Zonal flow) → BL is underestimated  
 NAO+ persistence is overestimated



Ensemble spread distribution at day 10 for forecasts initiated in: NAO+ (blue) blocking (red), NAO- (green) and AR (violet) regime



## Flow dependent verification:

- **Blocking is** the regime associated with the **least accurate forecasts**.
- Poor forecasts underestimate the persistence of blocking while overestimate the maintenance/transitions of/to zonal flow (NAO+)
- **The ensemble spread is a useful indicator of the forecast error.**
- **The spread of the forecasts initiated in NAO- is significantly smaller** than for the forecasts initiated in the other regimes. This is consistent with their higher skill.

# Predictability of Euro-Atlantic circulation regimes at extended range and its association to extreme events:

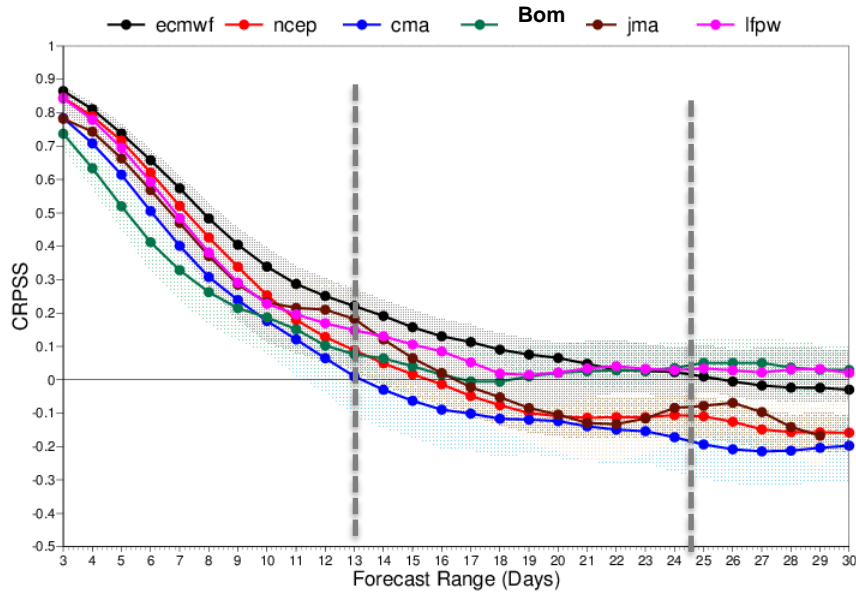
We are evaluating the predictive skill of the EA regimes using the S2S data base (Sub-seasonal to seasonal predictions WWRP/WCRP joint research project )

In particular we are interested in assessing the regime transitions ( climatological frequencies, loss of skill, physical processes associated with it)

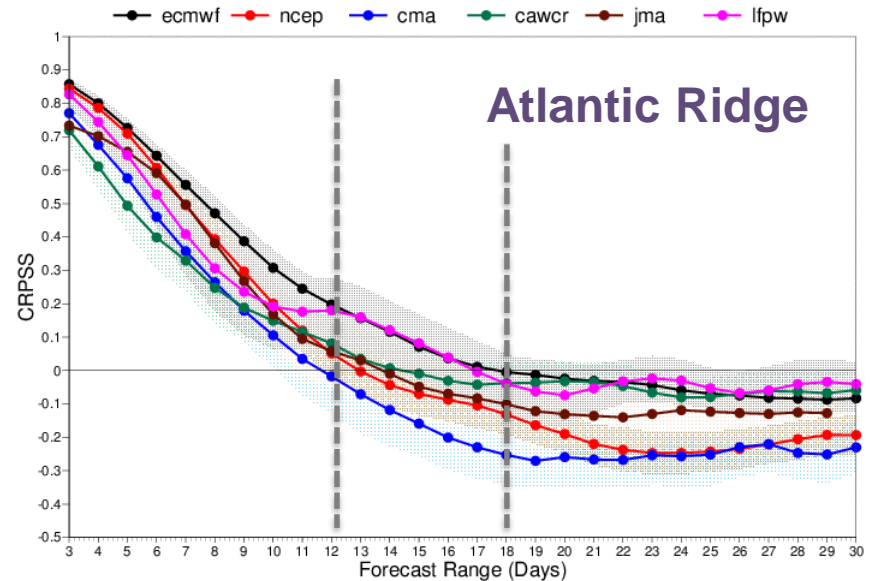
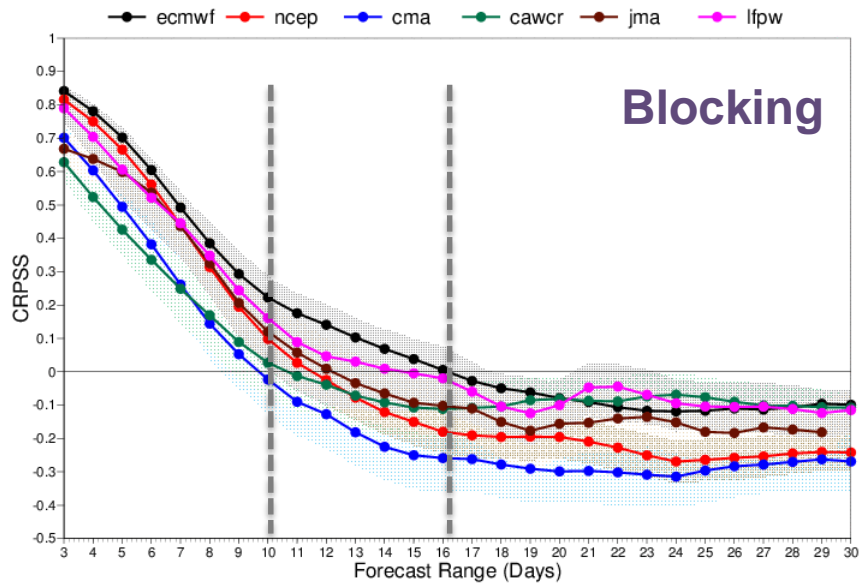
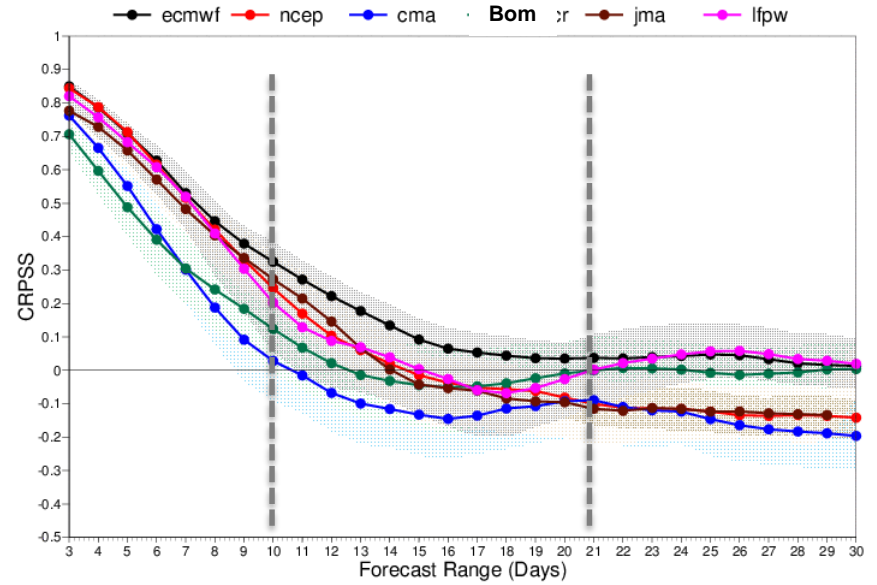
NAO- and BL are the flow patterns strongly associated with high impact temperature anomalies (heat waves in summer and cold spell in winter).

# Predicting skill associated with the Euro-Atlantic Regimes:

## NAO +

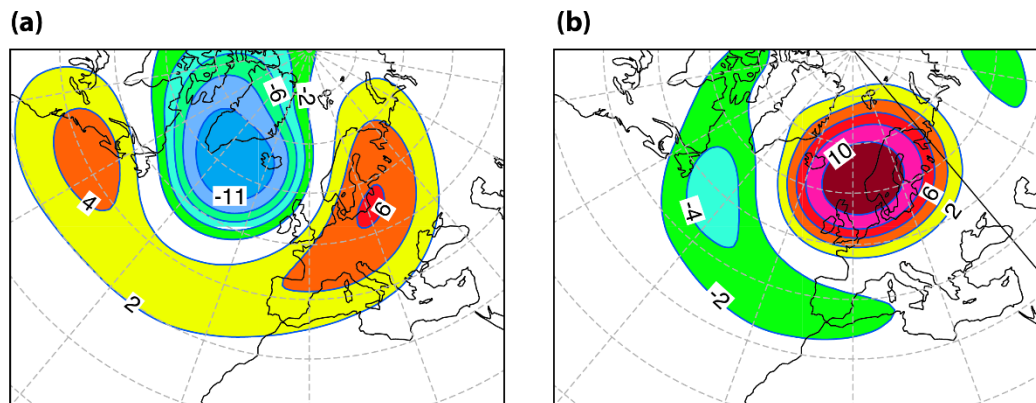


## NAO -

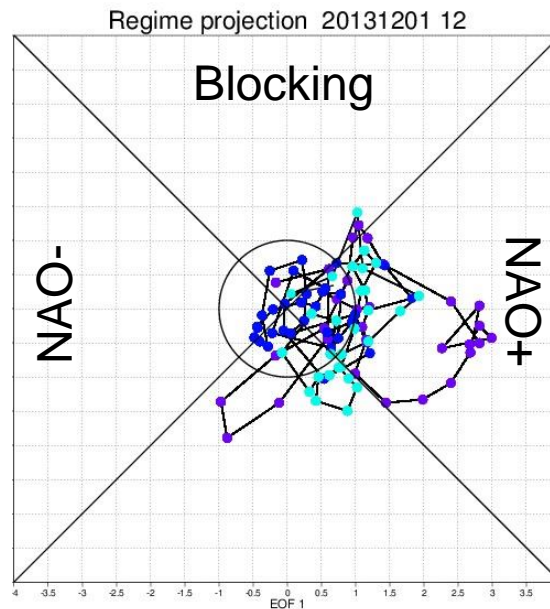
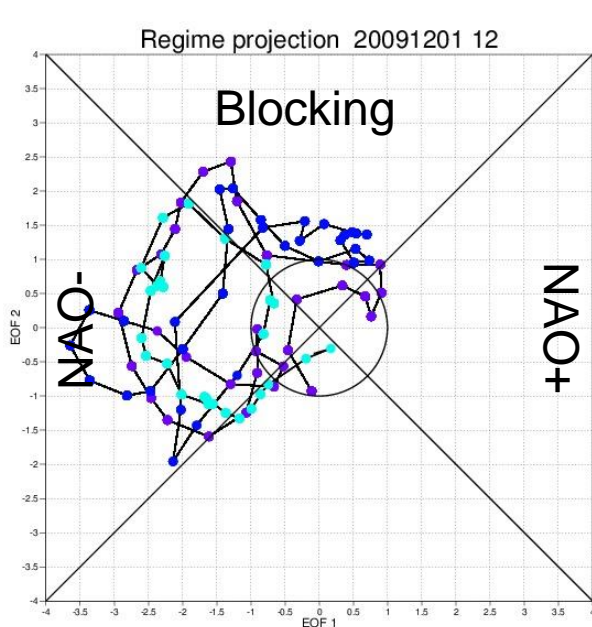


# How far in advance we predict changes in large scale flow leading to sever cold spell over Europe?

- $\pm$ EOF1 and  $\pm$ EOF2 represent quite well  $\pm$ NAO and BL
- Trajectories in phase space summarise regime evolution
- Unlike MJO, no preferred direction



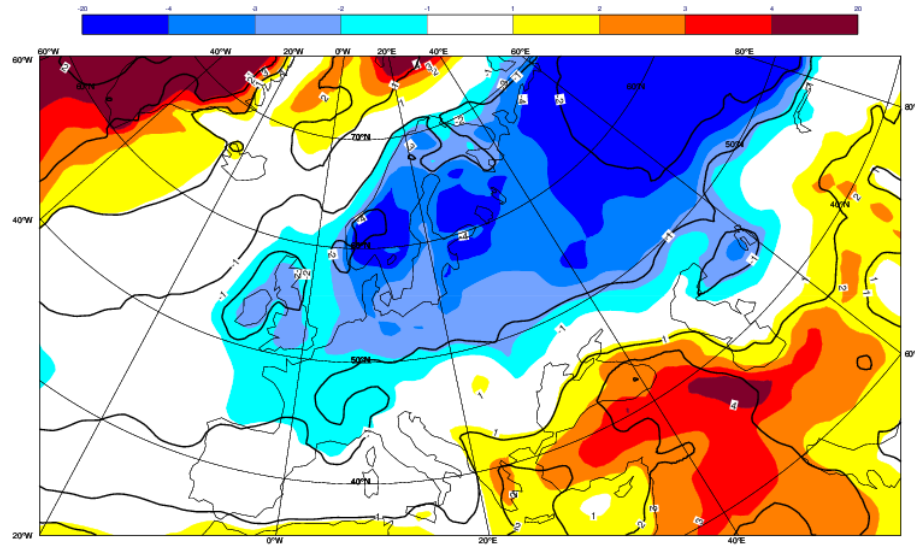
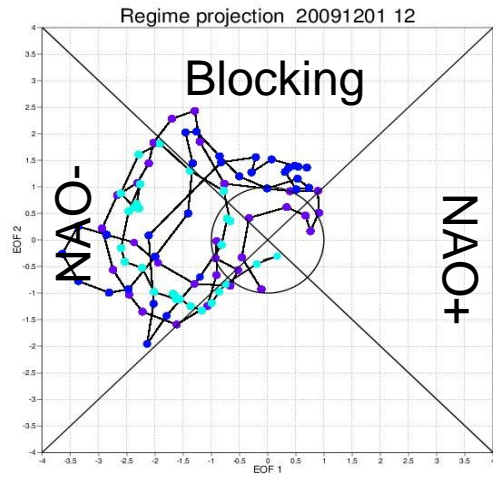
BL: record-breaking cold temperatures over Europe



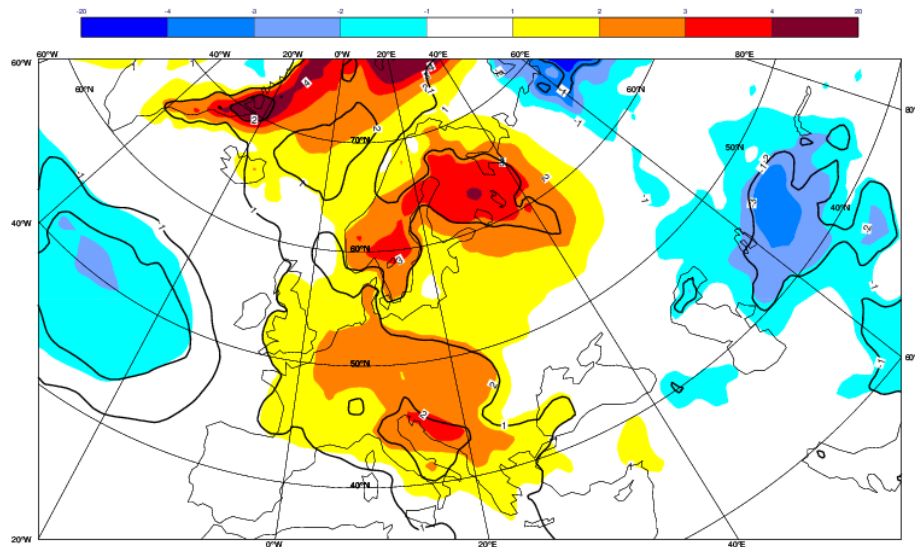
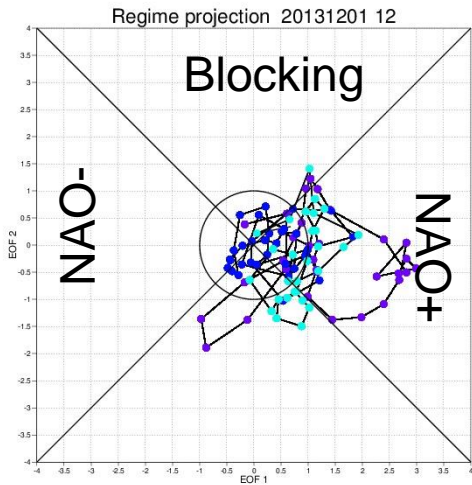
+NAO: exceptional storminess, but mild temperatures over Europe

# 2M Temp anomalies for DJF:

2009/2010



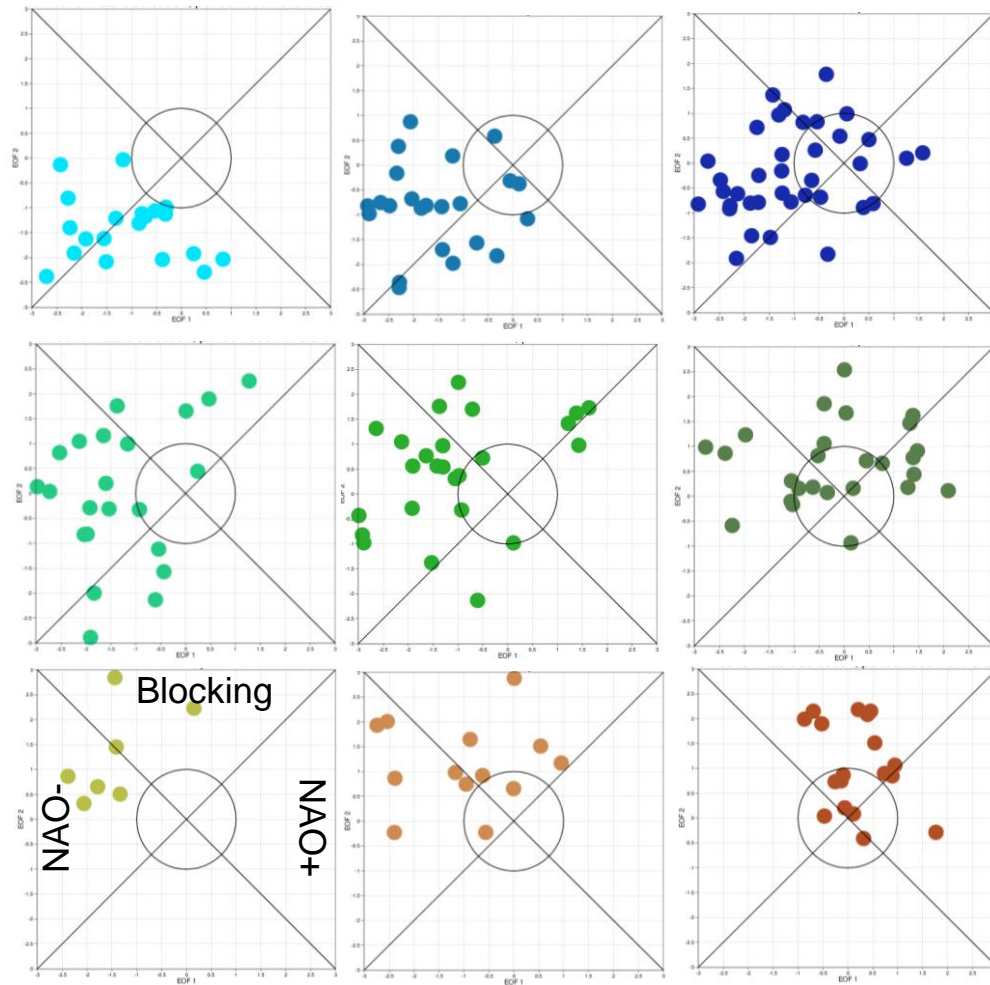
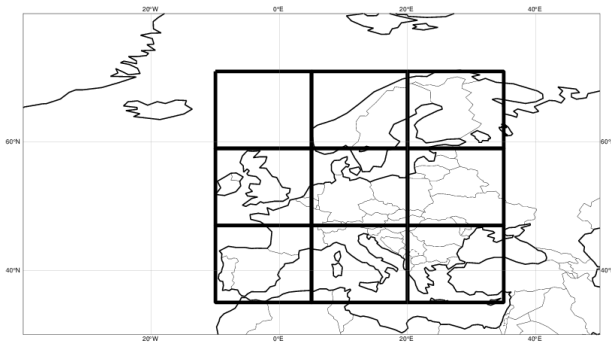
2013/2014



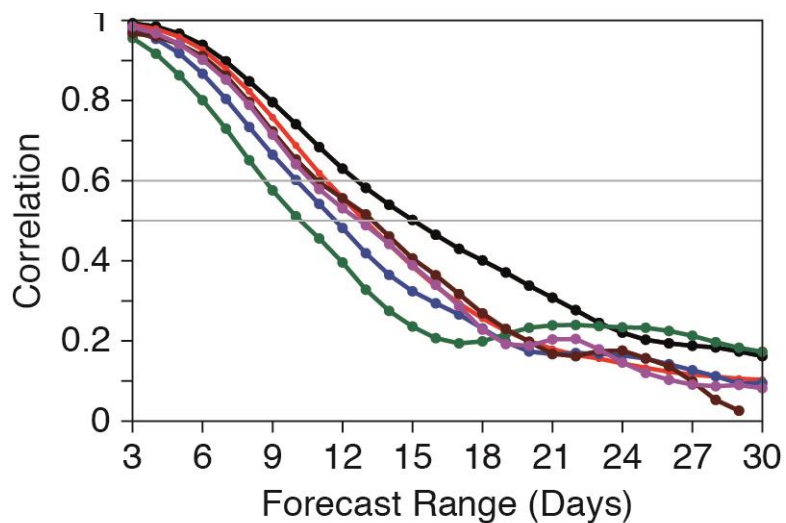


## Distribution of severe winter (NDJF) events in era-interim (1980-2015)

When for 60% grid points in each box the daily 2mt < 10<sup>th</sup> quantile of daily climate for at least 4 consecutive days

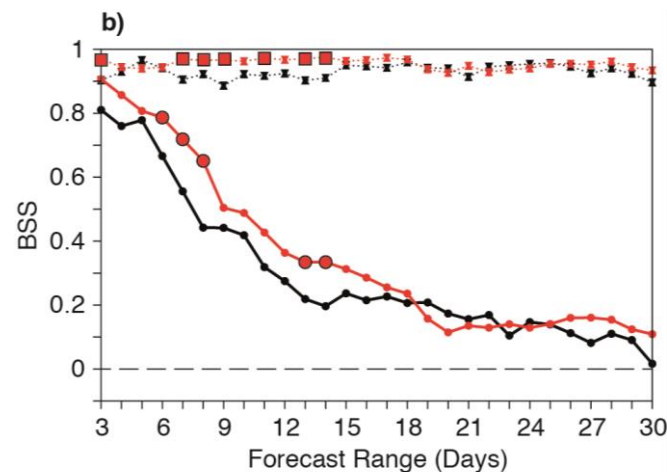


## Regime transitions:



—●— ECMWF    —●— NCEP    —●— CMA  
 —●— BoM      —●— JMA      —●— MetFr

## NAO- skill

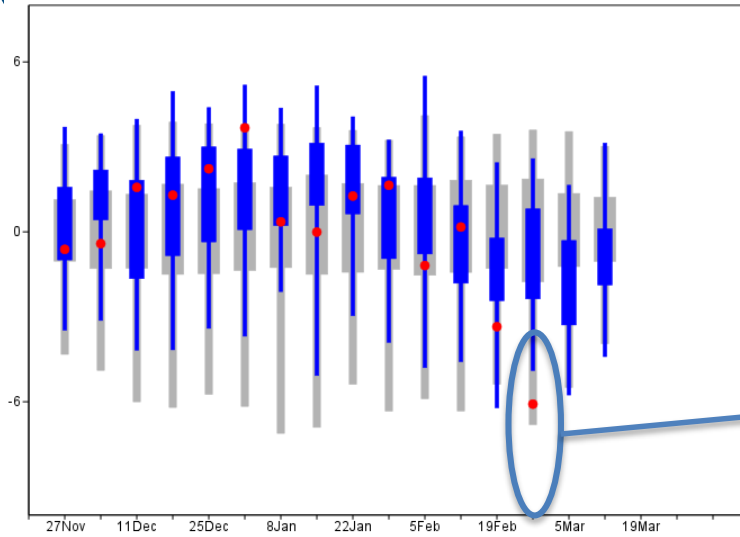


Small impact for NAO+ predictions  
 Significantly higher skill for NAO- forecasts with  
 and MJO in the i.c.

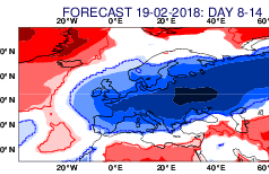
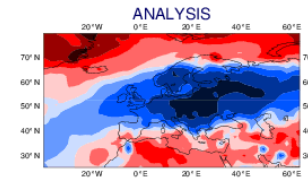
Ferranti, L. et al. 2018 *QJRMS*, 144, 1788–1802. doi:10.1002/qj.3341

# Severe cold spell end of February 2018:

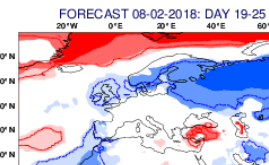
2mt over Europe  
weekly means anomalies at 19-25 days  
(3.5weeks)



26/2-4/3 2018



2 weeks ahead



3.5 weeks ahead

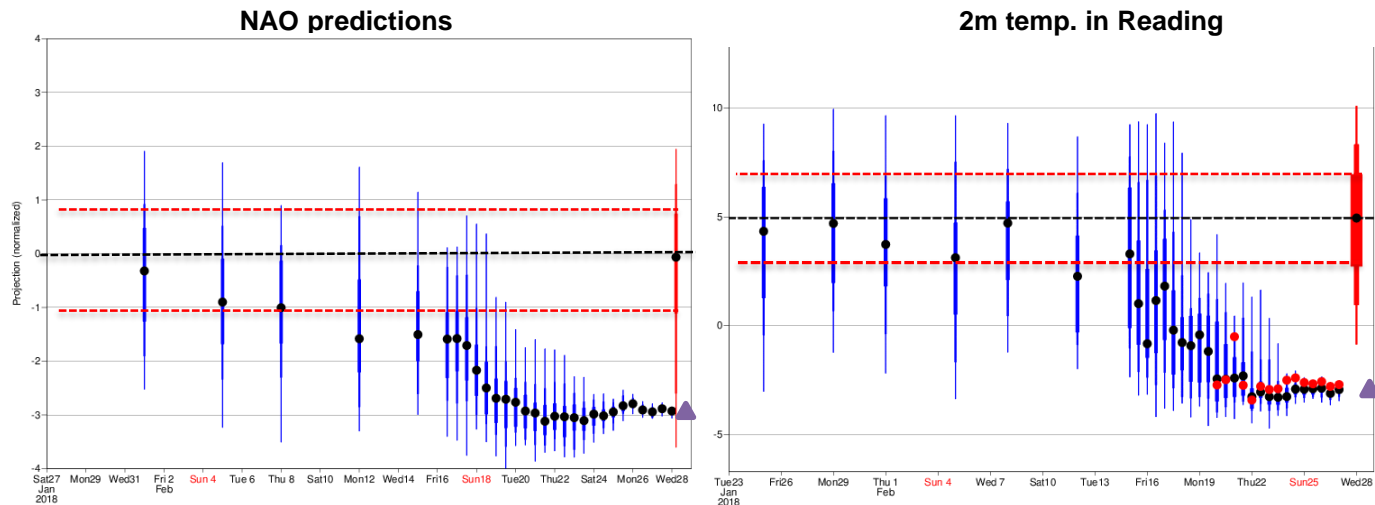




# Severe cold spell end of February 2018:

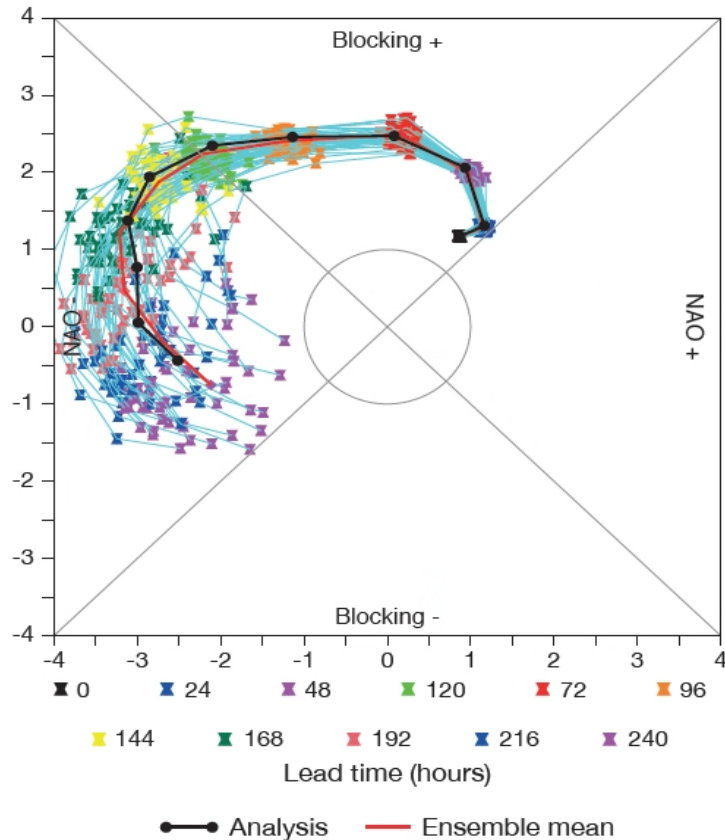
## How far in advance this cold event was predicted?

Predictions initialized at different time and verifying the 3-days mean (27 Feb to 1 March )



Circulation regimes, usually associated with global teleconnections, play an important role in the atmospheric predictability on sub-seasonal time scale. Regimes are associated with high impact events: cold spell in winter and heat waves in summer.

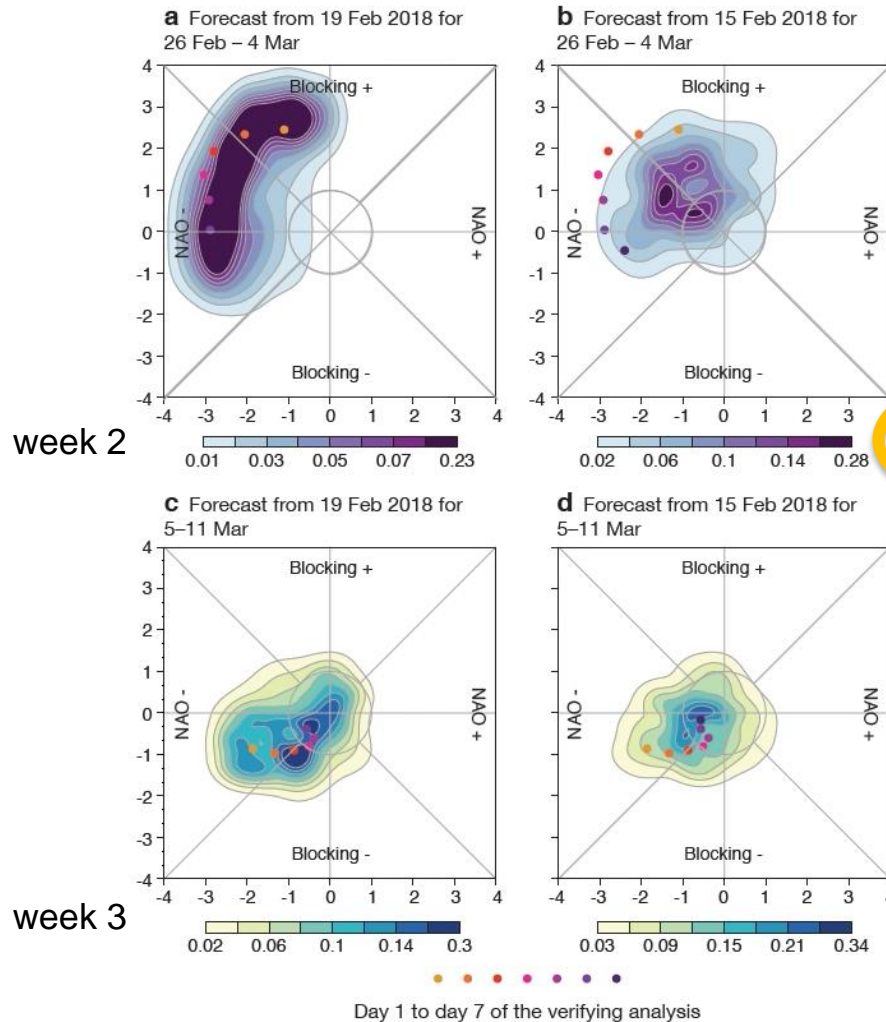
## Severe cold spell end of February 2018: Medium range forecast



Evolution of ENS forecast up to day 10 in the NAO–BLO diagram. The ENS forecast starts at 00 UTC on 22 February 2018.

See **ECMWF Newsletter 158** Winter 2018/19  
<https://www.ecmwf.int/en/publications>

# Severe cold spell end of February 2018: Extended range forecasts



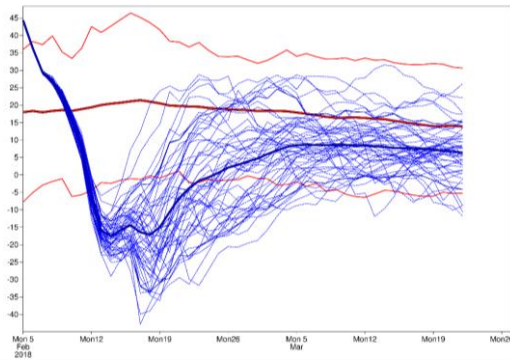
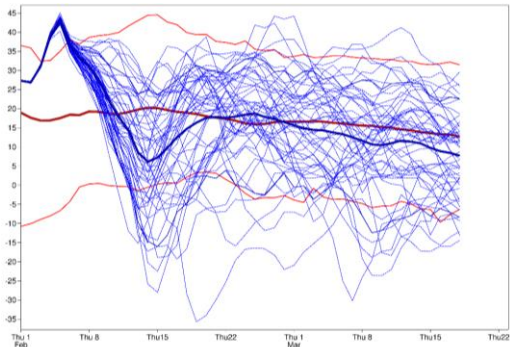
the accuracy of this forecast is link to MJO and SSW

**FIGURE 5** Probability density functions for (a) an ensemble forecast starting on 19 February 2018 for the week starting on 26 February, (b) an ensemble forecast starting on 15 February 2018 for the same week, (c) an ensemble forecast starting on 19 February 2018 for the week starting on 5 March and (d) an ensemble forecast starting on 15 February 2018 for the same week. Daily values of the verifying analysis are represented by dots.

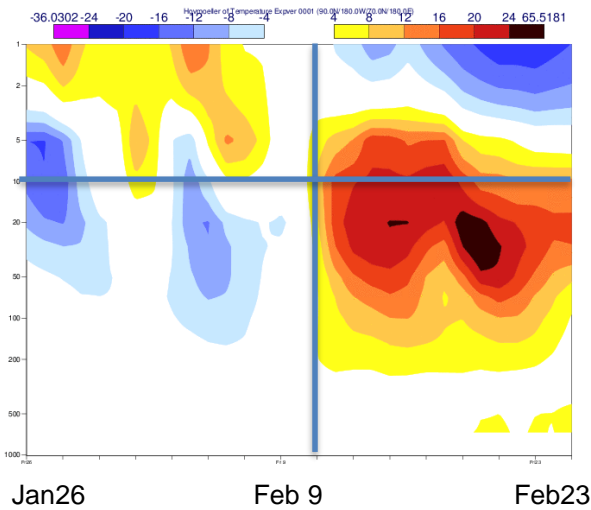
# Severe cold spell end of February 2018:

SSW:

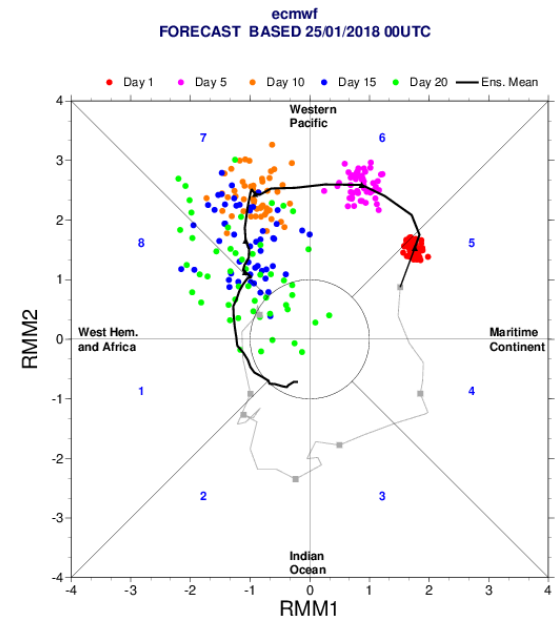
60N zonal mean zonal wind at 10hPa



11Feb SSW onset



MJO:



## Predictions of cold spells over Europe:

Reliable forecasts of NAO and blocking are instrumental for the extended range predictions of severe cold events over Europe.

S2S systems exhibit useful skill well beyond 10 days for NAO and Blocking predictions – strong potential for early warnings.

ECMWF forecasts, beyond 15 days, can provide reliable probabilities of cold temperatures associated with the NAO-.

Such skill can be enhanced by MJO activity (teleconnections) and SSW events.

Forecasting probabilities of cold spell associated with a blocking is a bigger challenge.

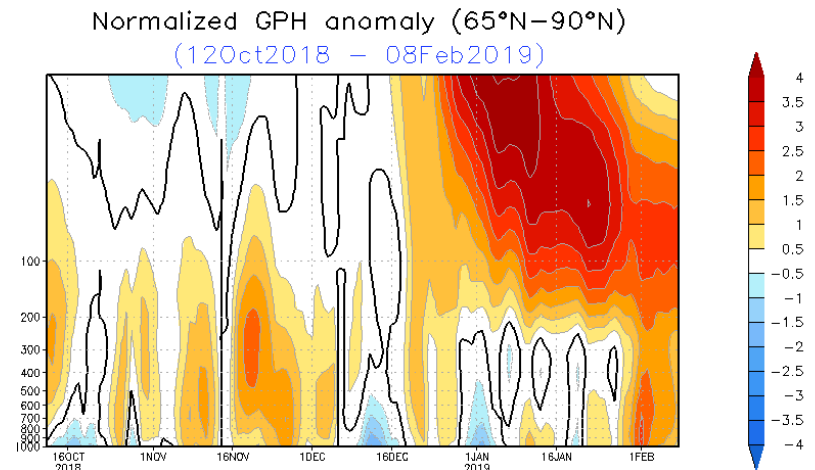
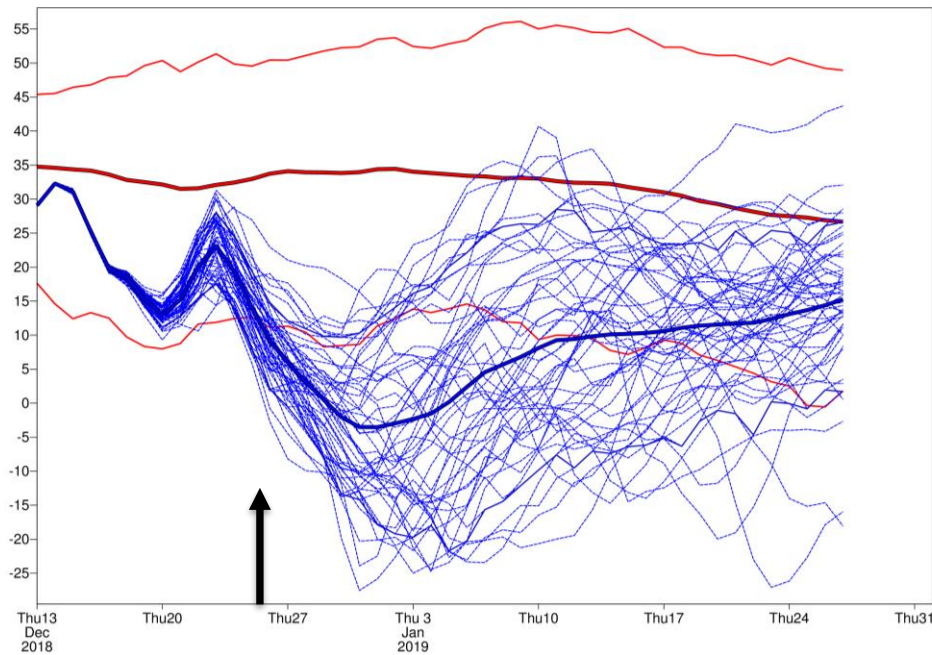
New test products :

<https://confluence.ecmwf.int/display/FCST/Test+products>

## The recent SSW event in Dec/Jan 2019

Forecast

Zonal mean zonal wind anomalies at 10hPa

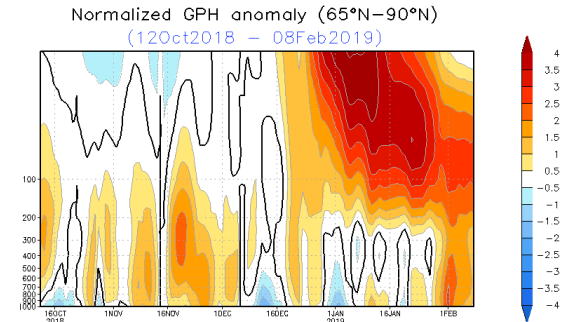
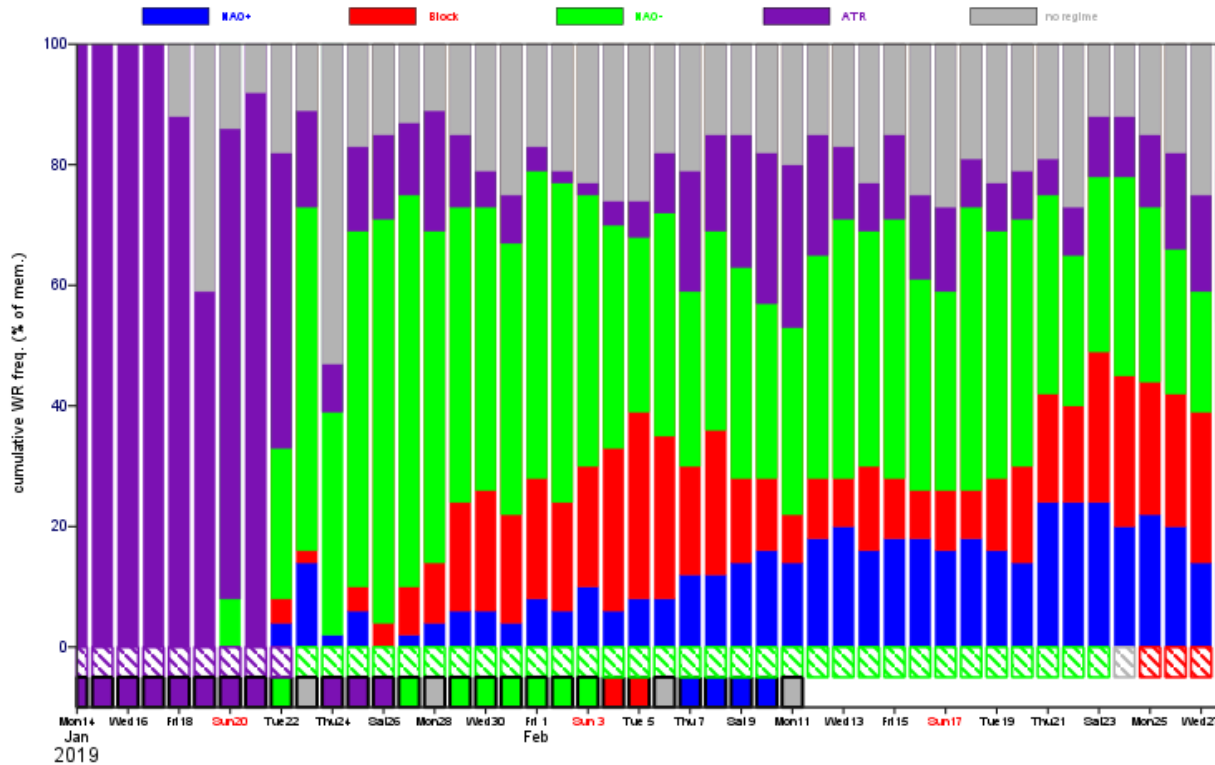


<https://www.cpc.ncep.noaa.gov>

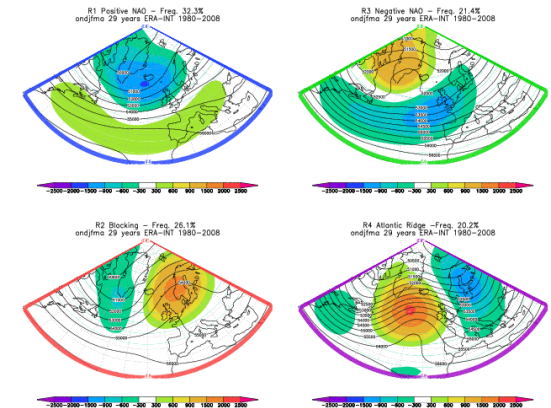
# New test products :

<https://confluence.ecmwf.int/display/FCST/Test+products>

Forecast Weather Regime frequency  
Ensemble Forecast: 20190114



<https://www.cpc.ncep.noaa.gov>



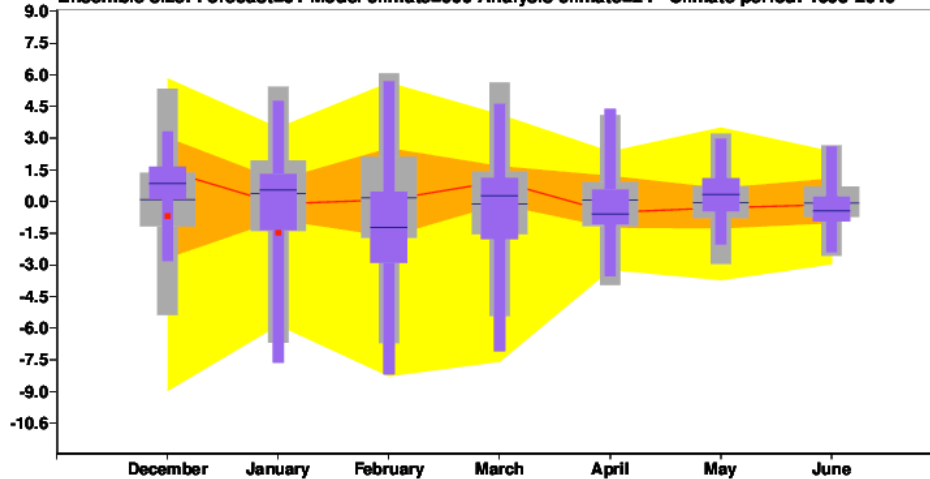


# Regime predictions at seasonal time scales: NAO seasonal forecasts

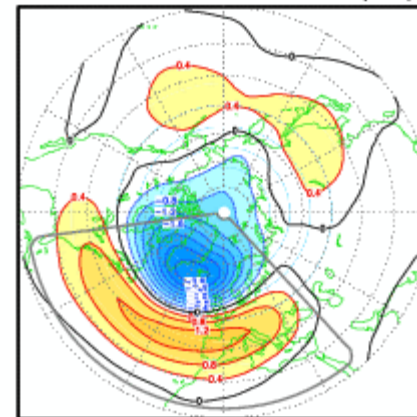
## North Atlantic Oscillation

Forecast Initial date: 20181201

Ensemble size: Forecast=51 Model climate=600 Analysis climate=24 Climate period: 1993-2016



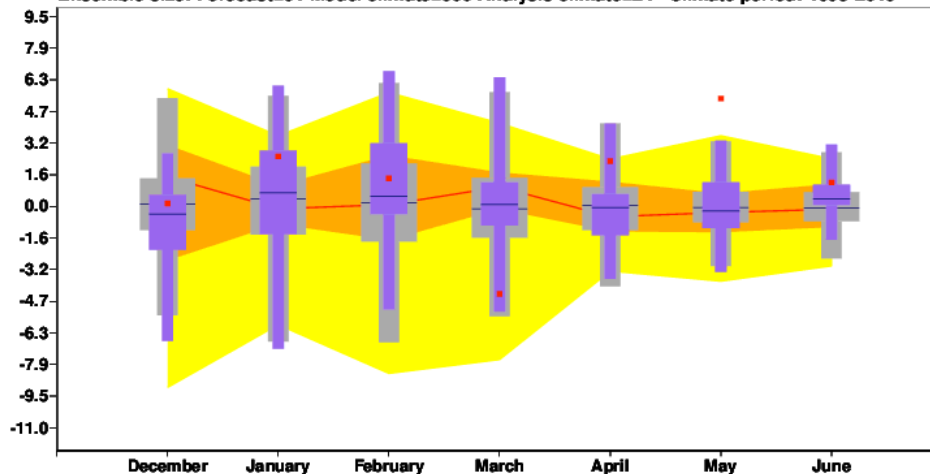
eof 1: North Atlantic Oscillation (NAO)



## North Atlantic Oscillation

Forecast Initial date: 20171201

Ensemble size: Forecast=51 Model climate=600 Analysis climate=24 Climate period: 1993-2016



**At seasonal scale  
skill for NAO is about 0.4**



# Summary

The skill in predicting E-A regimes, from medium to extended range, has been evaluated using several S2S reforecast systems - Some S2S models have the skill up to day10, other models show skill beyond the medium range. In terms of CRPSS ECMWF shows some skill for NAO-/NAO+ up to 20-23 days while for BL and AR skill drops to zero at about 16-17 days.

Forecast of transitions to/from regimes associated with high-impact temperature anomalies over Europe are evaluated using a 2-dim diagram.

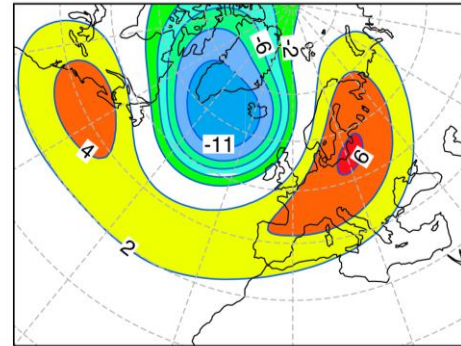
The forecast performance in predicting E-A regimes is monitored using the operational clustering products.

Questions?

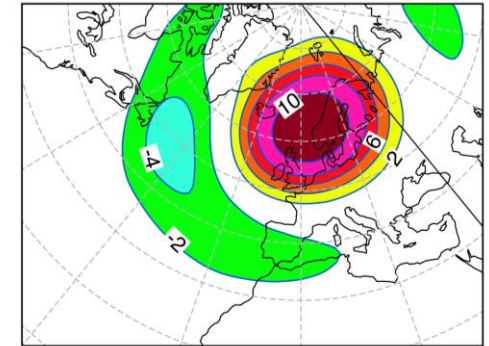


# NAO-BL diagrams

The ensemble evolution in the NAO-Blocking diagram :

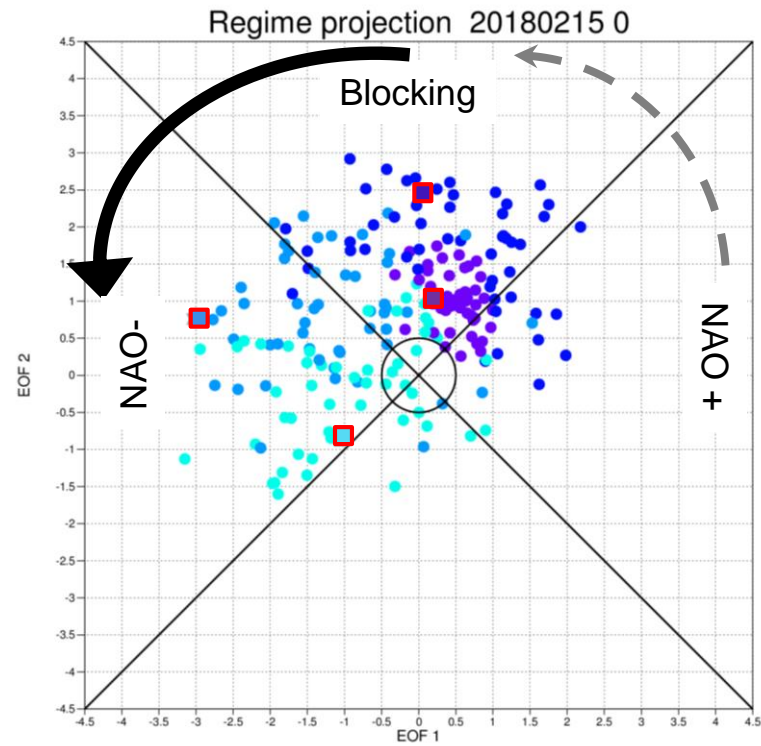
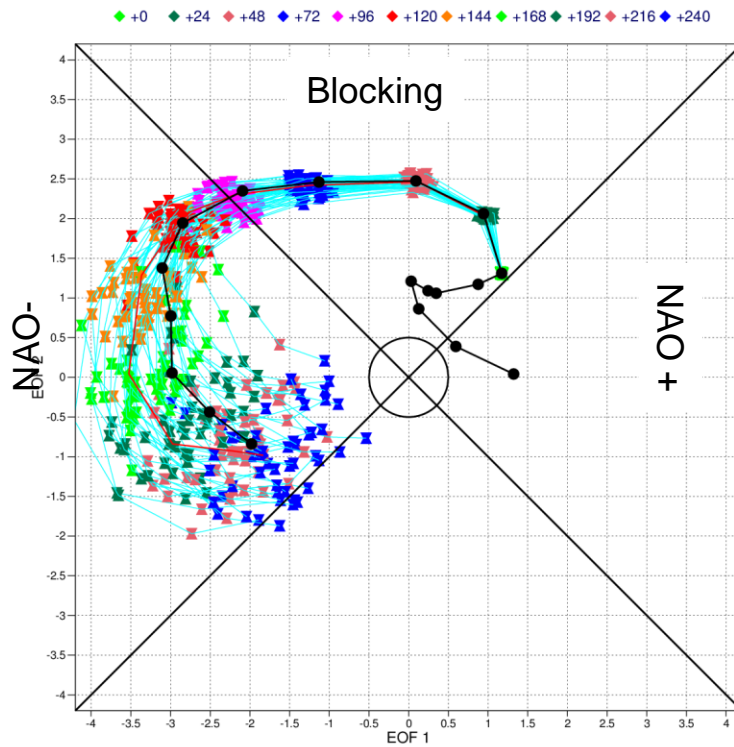


EOF1



EOF2

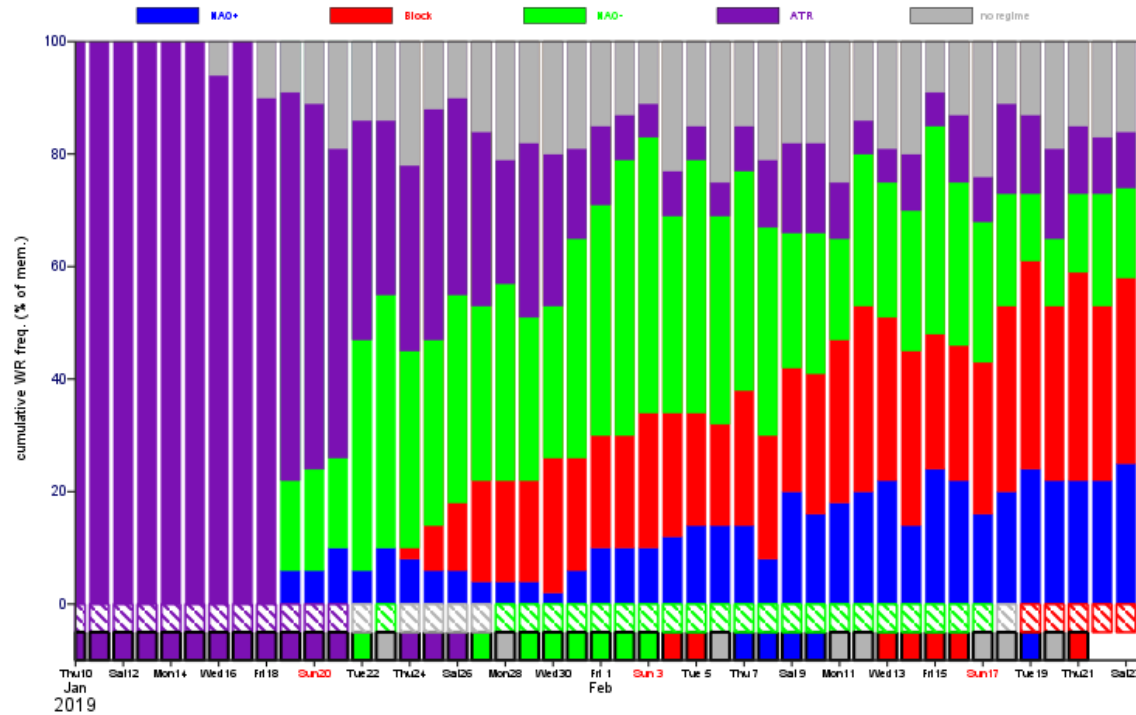
20170223  
Forecast



This winter we experience an NAO- event of massive amplitude that was predicted about 3 weeks in advance. The associated severe cold conditions were well represented by the weekly mean anomalies at 19-25.

The MJO and possibly the SSW might have played a role in enhancing predictability further analysis is needed.

## Weather Regime frequency Ensemble Forecast: 20190110



# Cluster scenario

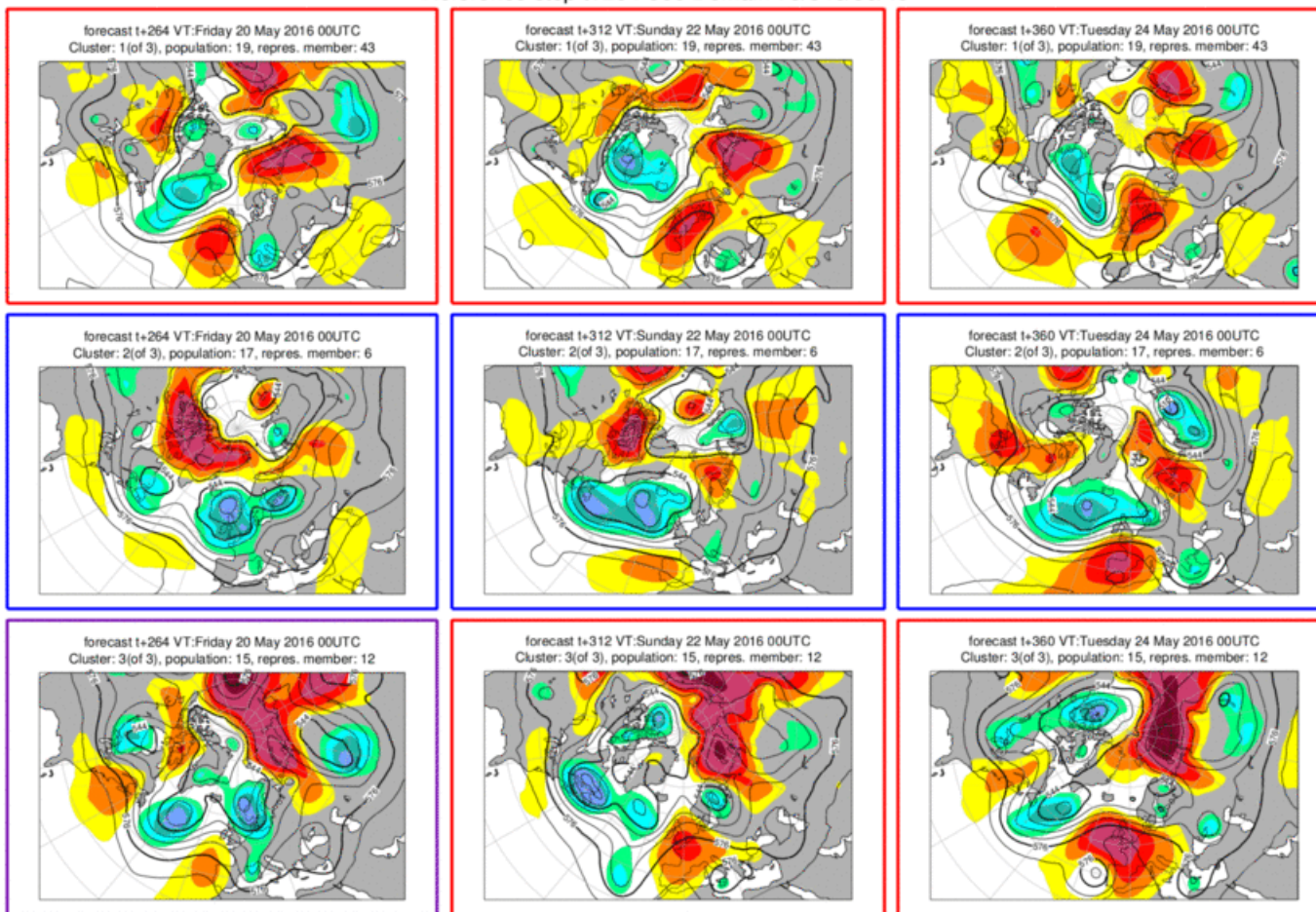
Base time

Parameter

Cluster

72-96 > 3-4 days  
 120-168 > 5-7 days  
 192-240 > 8-10 days  
**264-360 > 11-15 days**

Monday 9 May 2016 00UTC ECMWF EPS Cluster scenario - 500 hPa Geopotential  
 Reference step t+264-360 Domain 75/340/30/40

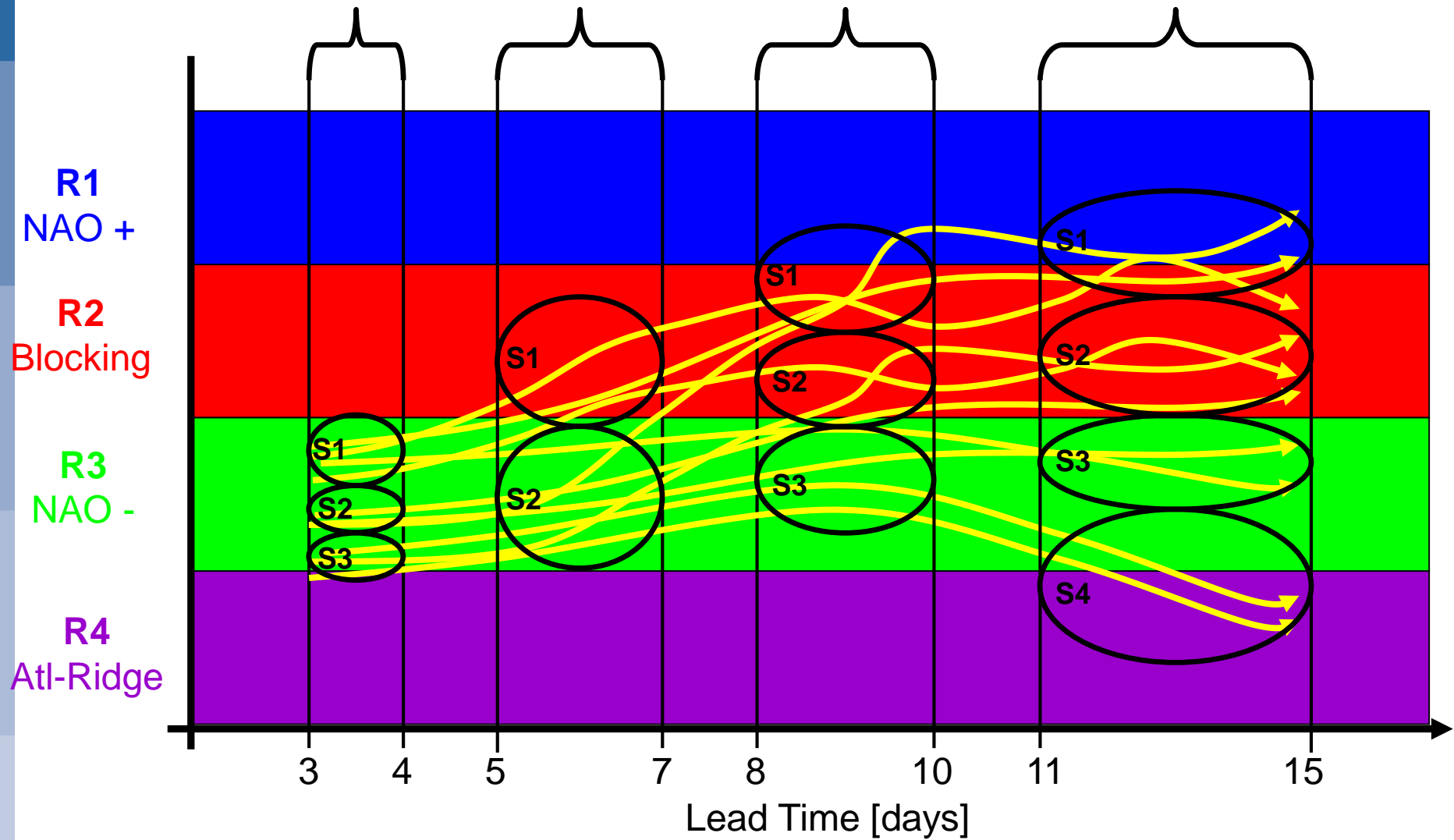


+11days

+13 days

+15 days

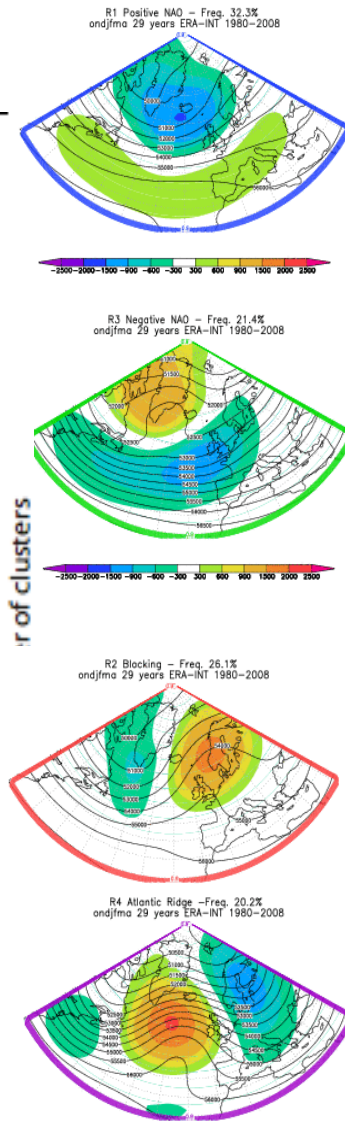
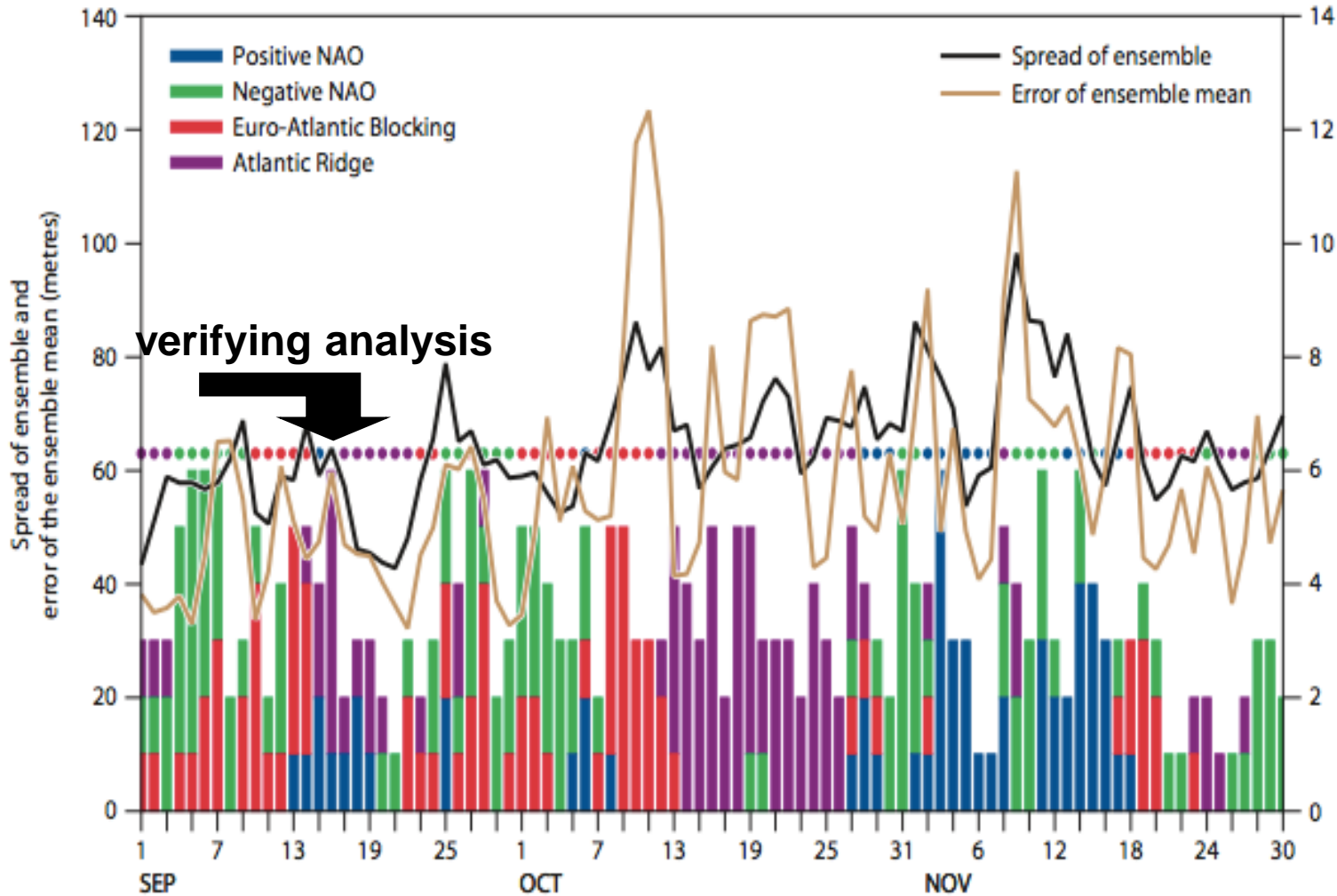
# Regimes & Scenarios

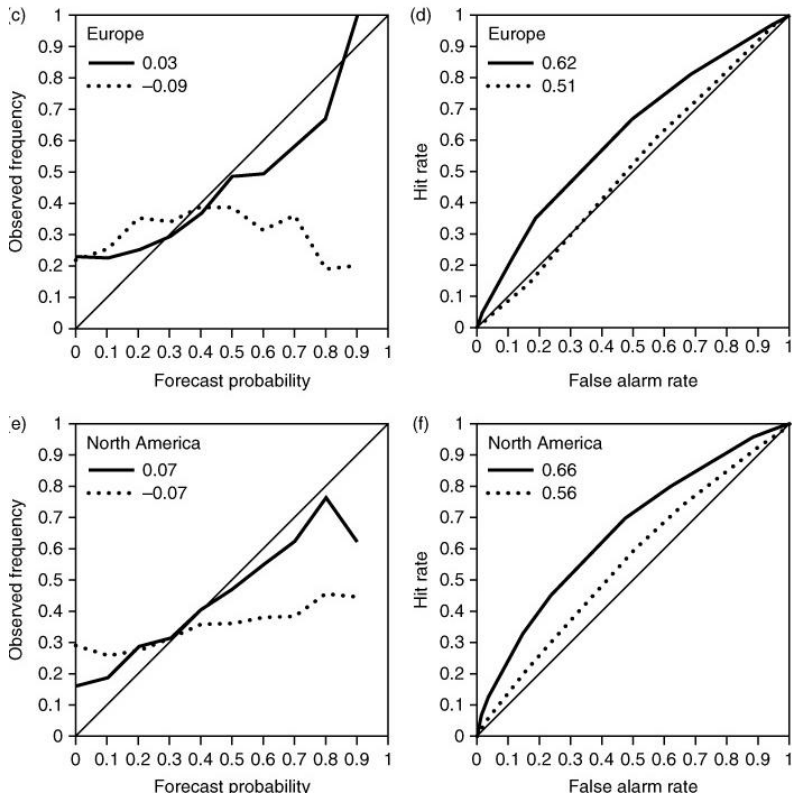




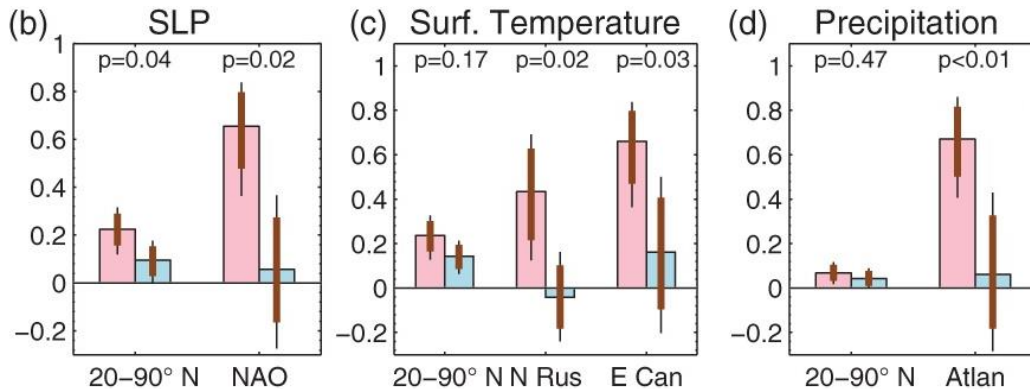
# Verification & spread

# Climatological regimes



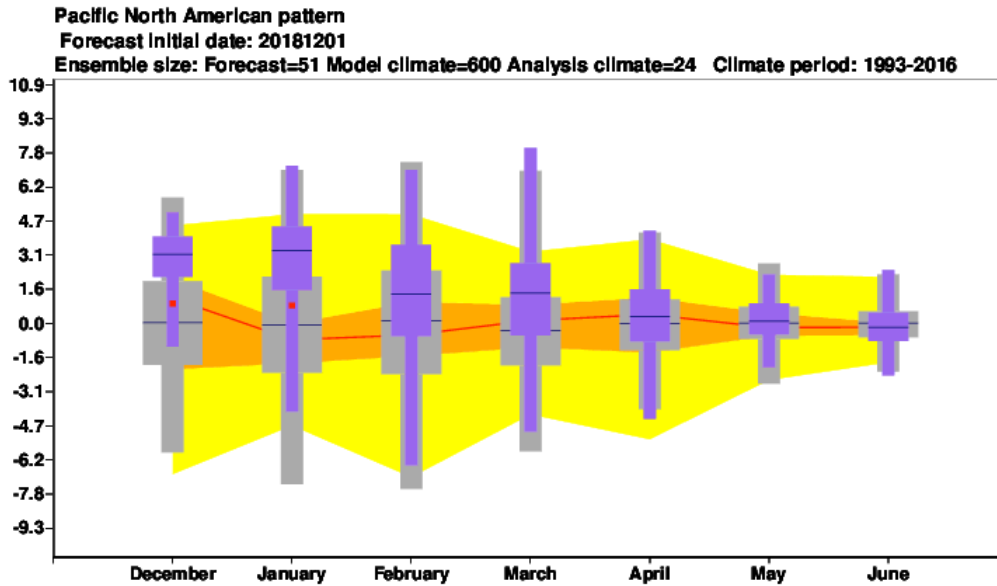


*From Vitart and Molteni 2010*

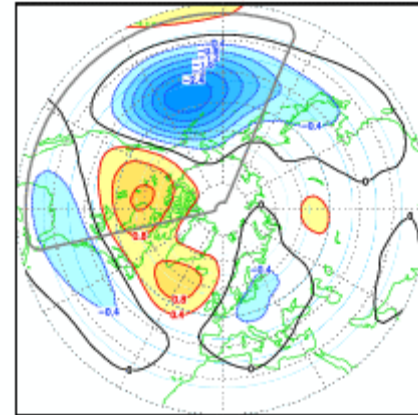


*From Tripathi et al. 2014*

# Regime predictions at seasonal time scales: Pacific North American pattern PNA



eof 1: Pacific/North-American (PNA)



**At seasonal scale  
skill for PNA is about 0.7**