

A standardized methodology for the validation of air quality forecast applications (F-MQO): Lessons learnt from its application across Europe

5 Lina Vitali¹, Kees Cuvelier^{2,a}, Antonio Piersanti¹, Alexandra Monteiro³, Mario Adani¹, Roberta Amorati⁴, Agnieszka Bartocha⁵, Alessandro D'Ausilio⁶, Paweł Durka⁷, Carla Gama³, Giulia Giovannini⁴, Stijn Janssen⁶, Tomasz Przybyła⁵, Michele Stortini⁴, Stijn Vranckx⁶, Philippe Thunis²

¹National Agency for New Technologies, Energy and Sustainable Economic Development (ENEA), Department for Sustainability, Bologna, Italy

10 ²European Commission - Joint Research Centre (JRC), Ispra, Italy

³CESAM, Department of Environment, University of Aveiro, Aveiro, Portugal

⁴Regional Agency for Prevention Environment and Energy (ARPAE) of the Emilia-Romagna Region, Bologna, Italy

⁵ATMOTERM, Opole, Poland

⁶Flemish Institute for Technological Research (VITO), Mol, Belgium

15 ⁷Institute of Environmental Protection (IEP) - National Research Institute, Warsaw, Poland

^aRetired with Active Senior Agreement

Correspondence to: Philippe Thunis (Philippe.THUNIS@ec.europa.eu)

Abstract. A standardized methodology for the validation of short-term air quality forecast applications was developed in the framework of FAIRMODE activities. The proposed approach, focusing on specific features to be checked when evaluating a forecasting application, investigates the model capability to detect sudden changes of pollutants concentrations levels, to predict threshold exceedances and to reproduce air quality indices. The proposed formulation relies on the definition of specific forecast Modelling Quality Objective and Performance Criteria, defining the minimum level of quality to be achieved by a forecasting application when it is used for policy purposes. The persistence model, which uses the most recent observed value as predicted value, is used as benchmark for the forecast evaluation. The validation protocol has been applied to several forecasting applications across Europe, using different modelling paradigms and covering a range of geographical contexts and spatial scales. The method is successful, with room for improvement, in highlighting shortcomings and strengths of forecasting applications. This provides a useful basis for using short-term air quality forecast as a supporting tool for correct information to citizens and regulators.

20
25

1 Introduction

30 Air pollution models play a key role in both enhancing the scientific understanding of atmospheric processes and supporting policy in adopting decisions aimed at reducing human exposure to air pollution. Current European Air Quality Directives

(AQD), 2008/50/EC (European Union, 2008) and 2004/107/EC (European Union, 2004), and even more the proposal of their revision (European Union, 2022), encourage the use of models in combination with monitoring in a wide range of applications. Indeed, models have the advantages of being cheaper than measurements and covering continuously and simultaneously large areas. Advances in the knowledge of atmospheric processes and the enhancement in computational technologies fostered the usage of 3-dimensional numerical models, the Chemical Transport Models (~~CTM~~), not only for air quality assessment (retrospective simulation of historical air quality scenarios in support of regulation and planning) but also for real-time air quality forecasting (~~AQF~~). Indeed, during last decades, air quality forecastingAQF systems based on ~~CTM~~ Chemical Transport Models have been rapidly developed and they are currently operational in many countries, providing early air quality warnings that allow policy makers and citizens to take measures in order to reduce human exposure to unhealthy levels of air pollution. On European scale, a real-time air quality forecasting system (Marécal et al., 2015) is operational since 2015 in the framework of the Copernicus Atmosphere Monitoring Service (CAMS) and currently includes eleven numerical air quality models, contributing to the CAMS Regional Ensemble production (<https://regional.atmosphere.copernicus.eu/>). Several review papers are available in literature, comprehensively describing current status and emerging challenges in real-time air quality forecasting (e.g. Kukkonen et al., 2012; Zhang et al., 2012; Baklanov et al., 2014; Ryan, 2016; Bai et al., 2018; Baklanov and Zhang, 2020; Sokhi et al., 2022), including air quality forecastingAQF system based on artificial intelligence (~~AI~~) methods (e.g. Cabaneros et al., 2019; Masood and Ahmad, 2021; Zhang et al., 2022).

A thorough assessment of model performances is fundamental to build confidence in models' capabilities and potentials and becomes imperative when model applications support policymaking. Moreover, performance evaluation is very important also for research purposes, since investigating models' strengths and limitations provides essential insights for planning new model developments.

The main goal of a model evaluation process is to prove that the performances are satisfactory for its intended use, in other words, that it is "fit for purpose" (e.g. Hanna and Chang, 2012; Dennis et al., 2010; Baklanov et al., 2014; Olesen, 1996). Indeed, to be able to determine whether a model application is "fit for purpose", its purpose should be stated at the outset. Since air quality models are used to perform various tasks (e.g. assessment, forecasting, planning), depending on the aim pursued, different evaluation strategies should be put into practice.

Several scientific studies have already proposed different evaluation protocols or suggested recommendations for good practices (e.g. Seigneur et al., 2000; Chang and Hanna, 2004; Borrego et al., 2008; Dennis et al., 2010; Baklanov et al., 2014; Emery et al., 2017). Models applied for regulatory air quality assessment are commonly evaluated by statistical analysis, examining how well they match the observations. From literature review, many statistical measures are used to quantify the different aspects of the agreement between simulations and observations. Indeed, no single metric is likely to reveal all aspects of model skills. So, the usage of several metrics, in concert, is generally recommended to support an in-depth assessment of performances. Zhang et al. (2012) provide an exhaustive collection of the most used metrics, ~~-. The list~~ includesing both traditional discrete statistical measures (e.g. Emery et al., 2017), quantifying the differences between

modelled and observed values, and categorical indices (e.g. Kang et al., 2005), describing the capability of the model application in predicting categorical answers (e.g. exceedances of limit values).

Ideally, a set of performances criteria should be given within a model evaluation exercise, stating if the model application skills can be considered adequate. As an example, Boylan and Russell (2006) and Chemel et al. (2010) proposed performance criteria and goals for mean fractional bias (*MFB*) and mean fractional error (*MFE*) concerning the validation of aerosol and ozone modelling applications, respectively. More in details, criteria define the acceptable accuracy level whereas goals specify the highest expected accuracy. Russell and Dennis (2000), citing Tesche et al. (1990), provided informal fitness criteria for urban photochemical modelling, according to some commonly used metrics (i.e. normalized bias, normalized gross error, unpaired peak prediction accuracy). Indeed, these recommendations are based on outcomes of performances skills from previous model studies. Specifically concerning [air quality forecasting AQP](#), in the framework of CAMS Regional Ensemble production, performances targets (Key Performance Indicators, KPI) are defined for the root mean square error (RMSE) in simulating ozone, nitrogen dioxide and aerosol, ~~and t~~ Their compliance is regularly reported within the Quarterly Evaluation and Quality Control Reports (<https://atmosphere.copernicus.eu/regional-services>).

Concerning both the definition of protocols for model evaluation and the proposal of performances criteria, an important contribution came in the last decades from the activities and the coordination efforts of the Forum for Air quality Modeling in Europe (FAIRMODE, <https://fairmode.jrc.ec.europa.eu/home/index>). FAIRMODE was launched in 2007 as a joint initiative of the European Environment Agency (EEA) and the European Commission Joint Research Centre. Its primary aim is to promote the exchange of good practices among air quality modellers and users and foster harmonization in the use of models by European Member States, with emphasis on model application under the European Air Quality Directives. In this context, one of the main activities of FAIRMODE has been the development of harmonized protocols for the validation and the benchmarking of modelling applications, ~~These protocols include~~ ~~ing~~ the definition of common standardized Modelling Quality Objectives (*MQO*) and Modelling Performance Criteria (*MPC*) to be fulfilled in order to ensure a sufficient level of quality of a given modelling application. More in details, an evaluation protocol was proposed for the evaluation of model applications for regulatory air quality assessment. The methodology (Thunis et al., 2012b; Pernigotti et al., 2013; Thunis et al., 2013; Janssen and Thunis, 2022) is based on the comparison of model-observation differences (namely, the root mean square error) with a quantity proportional to the measurement uncertainty. The rationale is that a model application can be considered ~~“fit for purpose”~~ ~~acceptable~~ if the model-measurement differences remain within a given proportion of the measurement uncertainty. The approach, ~~is~~ consolidated in the DELTA Tool software (Thunis et al. 2012a, <https://aqm.jrc.ec.europa.eu/Section/Assessment/Download>), ~~It~~ has reached a good level of maturity and ~~it~~ has been widely used and tested by model developers and users (Georgieva et al., 2015; Carnevale et al., 2015; Monteiro et al., 2018; Kushta et al., 2019). This approach focuses on applications related to air quality assessment, in the context of the AQD 2008/50/EC (European Union 2008), taking into account pollutants and metrics consistently with the AQD requirements.

Recently, FAIRMODE worked on developing and testing additional quality control indicators to be complied when evaluating a forecast application, extending the approach used for assessment applications. A scientific consensus was

100 reached, focusing on the model ability in the specific purpose of accurately predicting sudden changes and peaks in the pollutant concentration levels. The proposed methodology, based on the usage of the persistence model (e.g. Mittermaier, 2008) as a benchmark, is now publicly available for testing and application.

This paper describes this new standardized approach and is organized as follows. Sect. 2 illustrates the rationale and the main features of the developed methodology. Sect. 3 describes the setup of the forecasting simulations to which the methodology was applied, including information on the monitoring data used for the validation. Results are presented in Sect. 4, focusing on lessons learnt from the application of the proposed approach in different geographical contexts and spatial scales. Finally, some conclusions are drawn in Sect. 5 together with hints for further developments.

2 Methodology

The validation protocol proposed in this work is specific for forecasting evaluation. It is an extension of the consolidated and well documented methodology ~~The proposed methodology for forecasting evaluation comes on top of the consolidated evaluation protocol~~ fostered by FAIRMODE for the evaluation of model applications for regulatory air quality assessment. Therefore, it is recommended that the metrics suggested when evaluating forecasting applications are applied in addition to forecasting applications fulfil the standard assessment *MQO*, as defined in Janssen and Thunis (2022), ~~as well as the additional forecast objectives and criteria, as defined within the new specific protocol~~. This section describes ~~its~~ the main features of the proposed protocol, which focus ~~esing~~ ing on ~~some specific skills to be checked when evaluating a forecasting application, namely~~ the model capability to (1) detect sudden changes of concentrations levels (Sect. 2.1), ~~to~~ (2) predict threshold exceedances (Sect. 2.2) and ~~to~~ (3) reproduce air quality indices (Sect. 2.3). Note that the proposed approach is not exhaustive. It does not evaluate all relevant features of a forecast application and other analyses will be helpful to gain further insights into the behaviour, the strengths and the shortcomings of a forecast application.

120 The methodology, as currently implemented in the DELTA Tool software, supports the following pollutants and time averages: NO₂ daily maximum, O₃ daily maximum of 8-hour average, PM10 and PM2.5 daily mean.

2.1 Forecast Modelling Quality Objective (*MQO_f*) based on the comparison with the persistence model

Predicting the status of air quality is useful in order to prevent or reduce health impacts from acute episodes and to trigger short-term action plans. Therefore, it is of main interest to verify the forecast applications ability in getting the purpose of accurately reproducing sudden changes in the pollutant concentration levels. To account for this, within the proposed protocol, the main evaluation assessment of the “fitness for purpose” of a forecast application is based on the usage, as a benchmark, of the persistence model, that is by default not able to capture any changes in the concentration levels, since measurement data of the previous day are used as an estimate for the full forecast horizon. Indeed, the persistence approach is the simplest method for predicting the future behaviour, if no other information is available and is often used as a reference in verifying the performances of weather forecasts (e.g. Knaff and Landsea, 1997; Mittermaier, 2008).

Within the proposed forecasting evaluation protocol, the root mean square error of the forecast model is compared with the root mean square error of the persistence model. More in detail, a forecast Modelling Quality Indicator (MQI_f) is defined as the ratio between the two $RMSEs$, i.e.

$$MQI_f = \sqrt{\frac{\frac{1}{N} \sum_{i=1}^N (M_i - O_i)^2}{\frac{1}{N} \sum_{i=1}^N (P_i - O_i)^2}} \quad (1)$$

135 where M_i , P_i , O_i represent respectively the forecast, the persistence, and the measured values for day i , and N is the number of days included in the time series.

The persistence model uses the observations from the previous day as an estimate for all forecast days. As an example, we can consider a 3 day-forecast, providing today (day0), tomorrow (day1), and the day after tomorrow (day2) concentration values. If today is 5th February, persistence model uses data referring to yesterday (4th February) for all forecast data
 140 produced today. So, P_i refers to O_{i-1} for day0 (5th February), it refers to O_{i-2} for day1 (6th February) and it refers to O_{i-3} for day2 (7th February). More generally, the persistence model is related to the forecast horizon ($FH = 0, 1, 2$, etc.) as follows:

$$P_i = O_{i-1-FH} \pm U(O_{i-1-FH}) \quad (2)$$

where the measurement uncertainty U is also taken into account, consistently with the FAIRMODE approach. The methodology for estimating the measurement uncertainty as a function of the concentrations values is described in Janssen
 145 and Thunis (2022), where the parameters for its calculation for PM, NO₂ and O₃ are provided as well. It is important to note that we use as representative for the measurement uncertainty the 95th percentile highest value among all uncertainty values. For PM10 and PM2.5 the results of a JRC instrument inter-comparison (Pernigotti et al., 2013) have been used whereas a set of EU AIRBASE stations available for a series of meteorological years has been used for NO₂ and analytical relationships have been used for O₃. These 95th percentile uncertainties only include the instrumental error. More~~Some~~ details are
 150 provided in Appendix A.

The Ffulfilment of the forecast Modelling Quality Objective (MQO_f) is proposed ~~implies that~~ as a necessary but not sufficient quality test ~~minimum level of quality for policy purposes~~ isto be achieved by the forecasting application. The MQO_f is fulfilled when MQI_f is less than or equal to 1, indicating that the forecast model performs better (within the measurement uncertainty) than the persistence one, with respect to its capability of detecting sudden changes of concentrations levels.

155 Within the proposed protocol, two aspects are included in a single metric (MQI_f): (1) check how well the model prediction compares with measurements and (2) check whether the model prediction performs better than a given benchmark (here the persistence model).

The magnitude of the MQI_f score, since it is referenced to a benchmark, is dependent on the skill of the benchmark itself. To account for this, Aadditional Modelling Performance Indicators ($MPIs$) are ~~defined~~ proposed as part of the evaluation
 160 protocol, based on the mean fractional error (MFE), a normalized statistical indicator widely used in literature, defined as follows:

$$MFE = \frac{2}{N} \sum_{i=1}^N \frac{|M_i - O_i|}{(M_i + O_i)} \quad (3)$$

Based on this indicator, two different *MPIs* are defined and both included within the protocol: 1) $MPI_1 = MFE_f/MFE_p$ that compares the forecast model performances with the persistence model ones; 2) $MPI_2 = MFE_f/MF_U$ that evaluates forecast performances regardless of persistence aspects, using an acceptability threshold based on measurement uncertainty, where MF_U is the Mean Fractional Uncertainty, defined as follows:

$$MF_U = \frac{1}{N} \sum_{i=1}^N \frac{2U(O_i)}{O_i} \quad (4)$$

Using the uncertainty parameters provided in Table A1 in Appendix A, it turns out that $2U(O_i)/O_i$ shows larger values in the low concentration range and then tends towards a constant (0.5 for NO₂, 0.3 for O₃, 0.55 for PM10, 0.6 for PM2.5) at higher concentration values (Fig. A1, in Appendix A). So, the choice of MF_U as the acceptability threshold is consistent with performances criteria and goals defined in literature for PM (Boylan and Russell, 2006) and O₃ (Chemel et al., 2010) and it has the advantage that it does not introduce any additional free parameters and it can be applied to all pollutants for which uncertainty parameters are set. For both *MPIs*, Modelling Performance Criteria (*MPC*) are ~~defined~~proposed, being fulfilled when *MPIs* are less or equal to 1.

2.2 Assessment of modelling application capability in predicting Threshold Exceedances

When a forecasting system is used for policy purposes, it is of main interest to verify the skill in predicting categorical answers (yes/no) in relation to exceedances of specific threshold levels, e.g. the limit values set by the current European legislation (European Union, 2008).

To account for this, the most commonly used threshold indicators (as defined in Table 1) are included in the proposed validation approach, based on the 2x2 contingency table (Table B1, in Appendix B) representing the joint distribution of categorical events (below/above the threshold value) predicted by the model and observed by the measurements. Namely, GA_+ represents the number of correctly forecasted exceedances, GA_- represents the number of correctly forecasted non-exceedances, FA (False Alarms) represents the number of forecasted exceedances that were not observed, and MA (Missed Alarms) represents the number of observed exceedances that were not forecasted.

All metrics included are listed in Table 1, ranging from 0 to 1, being 1 the optimal value.

Table 1. Categorical metrics included in the validation protocol.

Metrics	Mathematical Expression
Accuracy	$ACC = \frac{GA_+ + GA_-}{FA + GA_+ + GA_- + MA}$

Success Ratio	$SR = \frac{GA_+}{FA + GA_+}$
Probability of Detection	$PD = \frac{GA_+}{GA_+ + MA}$
FBias score	$FB = \frac{FA + GA_+}{GA_+ + MA}$
Threat score	$TS = \frac{GA_+}{FA + GA_+ + MA}$
Gilbert Skill score	$GSS = \frac{GA_+ - H}{FA + GA_+ + MA - H}$ <i>with</i> $H = \frac{(GA_+ + MA)(FA + GA_+)}{FA + GA_+ + GA_- + MA}$

2.3 Assessment of modelling application capability in predicting Air Quality Indices

190 One of the main objectives of a forecasting system is to provide citizens with simple information about local air quality and its potential impact on their health, with special regard to the sensitive and vulnerable groups (i.e., the very young or old, asthmatics, etc.). Air Quality Indices (AQI) are designed to provide information on the potential effects of the different pollutants on people's health by means of a classification of concentrations values in terms of qualitative categories.

The AQI outcome is commonly provided by operational forecasting systems, therefore its assessment has been included in the proposed validation approach, by means of a simple multiple thresholds assessment. More in detail, the number of days 195 predicted by the forecast model in each category is compared with the corresponding number of measured days.

Of course, the performance assessment depends on the chosen classification table. In the current approach, several AQI tables are available, namely EEA (<https://www.eea.europa.eu/themes/air/air-quality-index/index>), United Kingdom (<https://uk-air.defra.gov.uk/air-pollution/daqi>; <https://uk-air.defra.gov.uk/air-pollution/daqi?view=more-info>), and USEPA (U.S. Environmental Protection Agency, <https://www.airnow.gov/aqi/aqi-basics/>, Eder et al. 2010) classification tables.

200 3 Forecasting applications: models, setup and monitoring data for validation

The proposed methodology was applied across Europe to evaluate the performances of several forecasting applications. This paper focuses on lessons learnt by the validation of five forecasting applications, based on various methods (both in terms of chemical transport models and statistical approaches) and covering different geographical contexts and spatial scales, from

very local to European scale. The key features of the forecast applications are summarized in Table 2. Some more details are
205 provided for each of them in the following, along with information on the monitoring data used for the validation.

Table 2. Main features of the forecast applications.

Forecast application Acronym	Operated by	Modelling System	Modelling approach	Time Period	Horizontal Domain & Resolution	Meteo	Emissions	Boundary Conditions	Data Assimilation
FA1	ENEA	MINNI	<u>Dispersion</u> <u>Chemical</u> <u>Transport</u> Model	Year-long Simulation (2018)	Europe (25°W-45°E, 30°N-72°N) Resolution: 0.1°	IFS	CAMS REG (v5.1)	C-IFS	NO
FA2	CESAM	WRF-CHIMERE	<u>Dispersion</u> <u>Chemical</u> <u>Transport</u> Model	Year-long Simulation (2021)	Portugal (10.3°W-5.7°W, 36.4°N-42.6°N) Resolution: 0.05°	NCEP/GFS	EMEP/CEIP	GOCART for dust, LMDz-INCA for gaseous and other aerosol species	NO
FA3	VITO	OPAQ	Neural Networks	Year-long simulation (2022)	Ireland (10.5°W-5.9°W, 51.4°N-55.4°N) Resolution: 3 km	ECMWF	Not applicable	Not applicable	NO
FA4	ARPAE	NINFA	<u>Dispersion</u> <u>Chemical</u> <u>Transport</u> Model	Year-long Simulation (2021)	Prepair -PREPAIR domain (6.25°E-16.75° E, 43.1°N-47.35°N) Resolution: 0.07° x 0.05°	COSMO	ISPRA, <u>EMEP/CEIP,</u> PREPAIR	kAIROS	NO
FA5	ATMOTERM	CALPUFF	Dispersion Model	July 2020 - September 2022	Variable spatial grid-size covering Kosovo (1 km/0.5km) Pristine (200m/50m)	WRF	Kosovo emission inventory	CAMS ENSEMBLE	YES

3.1 MINNI simulation over Europe (FA1)

The first forecast application (FA1) was operated by ENEA applying the MINNI Atmospheric Modelling System (Mircea et al., 2014; D'Elia et al., 2021) on a European domain at 0.1° horizontal spatial resolution. ~~FA1 is a~~ year-long simulation, referring to 2018, ~~was carried out using CAMS Regional Emission Inventory (CAMS-REG, version v51, <https://eccad3.sedoo.fr/catalogue>), daily biomass burning emissions from Global Fire Assimilation System (GFAS, <https://confluence.ecmwf.int/display/CKB/CAMS+global+biomass+burning+emissions+based+on+fire+radiative+power+%28GFAS%29%3A+data+documentation>), Integrated Forecasting System meteorological fields (IFS, <https://www.ecmwf.int/en/publications/ifs+documentation>) and boundary conditions from the Chemical Integrated Forecasting System (C-IFS; <https://www.ecmwf.int/en/research/modelling+and+prediction/atmospheric+composition>).~~ MINNI, which is operationally providing air quality predictions over an Italian domain since 2017 (Adani et al., 2020, 2022), was recently added to the ensemble of the eleven models contributing to the CAMS Regional Ensemble production. ~~As~~ FA1 was carried out during a preliminary benchmark phase, ~~planned for evaluating model performances and for setting up the operational chain, it was carried out~~ using CAMS input and setup, but it is not an official CAMS product. Since no data assimilation was applied within FA1, all available data measured at European background monitoring stations and collected by EEA (E1a at <https://discomap.eea.europa.eu/map/fme/AirQualityExport.htm>) were considered for the validation.

3.2 WRF-CHIMERE simulation over Portugal (FA2)

In Portugal, an air quality modelling system based on the WRF version 3 (Skamarock et al., 2008) and the CHIMERE chemical transport model v2016a1 (Menuet et al., 2013; Mailler et al., 2017) is being used for forecasting purposes at daily basis since 2007 (Monteiro et al., 2005, 2007a, b). The modelling setup comprises three nested domains covering part of the North Africa and Europe, with horizontal resolutions of 125×~~125~~ km², 25×~~25~~ km² and 5×~~5~~ km² for the innermost domain covering Portugal. At the boundaries of the outermost domain, the outputs from LMDz-INCA (Szopa et al., 2009) are used for all gaseous and aerosol species, and dust from the GOCART model (Ginoux et al., 2001). ~~For the nested domains, the boundary conditions are updated every hour using the outputs from the coarse CHIMERE simulation. The initial conditions are defined by the 24 h forecast from the previous day model run. In addition to the meteorological fields (from the WRF Model) and the chemical boundary conditions, another important input includes the primary pollutant emissions.~~ The main human activities emissions (traffic, industries and agriculture, among others) are derived based on data from the annual EMEP/CEIP emission database (available at <https://www.ceip.at/webdab-emission-database/>), following a procedure of spatial and temporal downscaling. Biogenic emissions are computed online using the MEGAN model (Guenther et al., 2006), while dust emission fluxes are calculated using the dust production model proposed by Alfaro and Gomes (2001).

Data from the Air Quality National Monitoring network (<https://qualar.apambiente.pt>) is used every year to assess the performance of this forecasting modelling system, usually evaluated at annual basis. This comprehends a group of more than 40 background monitoring stations, classified as urban, suburban and rural environment, according to the classification settled by European legislation.

245 3.3 OPAQ simulation over Ireland (FA3)

The OPAQ (Hooyberghs et al., 2005; Agarwal et al., 2020) statistical forecast system has been configured and applied to forecast pollution levels in Ireland by the Irish EPA and VITO. During the configuration stage neural networks are trained at station level with historical observations, ECMWF-ERA5 reanalysis meteorological data and the CAMS air quality forecasts. ~~For each monitoring station, different neural network models are tested and validated. Both feed forward (FF) and recurrent neural networks are applied. The model showing the best performances will be operationally used to inform the public on the expected evolution in pollution levels for a forecast horizon available up to three days on an hourly basis.~~ The forecasts at station level are interpolated to forecast maps for the whole country using the detrended kriging model RIO (Janssen et al., 2008; Rahman et al., 2023) which is part of the OPAQ system. ~~The interpolation model has been configured for Ireland at 3 by 3 km². A spatial driver for the detrending and retrending step is constructed and optimized relying on CORINE 2012 land cover data, population density, altitude information and CAMS reanalysis annual average concentration maps.~~

In this study, we present the historical validation results of a feed-forward neural network model that uses 2-metre temperature, vertical and horizontal wind velocity component, CAMS PM10 forecasts, and PM10 observations. More than two years of data are used to configure the OPAQ model. Data from October 2019 to June 2022 are used for training. The model is validated on the data for July to December 2022. The testing holdout sample, used to avoid over fitting, covers a timespan of three months from June to September 2019. The model was optimized using the Adamax algorithm (Kingma and Ba, 2014) with 4 hidden layers and 200 units per layer, the activation function uses sigmoid functions while the mean squared error is used as loss function.

3.4 NINFA simulation over Po Valley and Slovenia (FA4)

FA4 was operated by ARPAE applying NINFA, his operational Air Quality Model Chain over Po Valley and Slovenia in the framework of Life Ip ~~Prepair~~ PREPAIR project (<https://www.lifeprepare.eu/>; Raffaelli et al., 2020), ~~with horizontal resolution of 5 km.~~ The model suite includes a Chemical Transport Model, a meteorological model and an emissions pre-processing tool. The chemical transport model is CHIMERE, v2017r3. ~~Starting from the eEmission data for cover~~ the Po Valley (Marongiu et al., 2022), Slovenia and the other regions/countries present in the model domain (http://www.lifeprepare.eu/wp-content/uploads/2017/06/Emissions-dataset_final-report.pdf), ~~the emissions are prescribed to the grid model by using specific proxy variables for each emission activity at SNAP3 level (i.e. road network for traffic emission, population and urban fabric for domestic heating, and so on).~~ The meteorological hourly input is provided by

COSMO (<http://www.cosmo-model.org>; Baldauf et al., 2011; Doms and Baldauf, 2018), ~~the national model used by the National Civil Protection Department. COSMO is a non-hydrostatic, limited-area atmospheric prediction model, based on the primitive thermo-hydrodynamical equations describing compressible flow in a moist atmosphere, with a variety of physical processes taken into account by dry and moist parameterization schemes.~~ The boundary conditions are provided by kAIROS (Stortini et al., 2020) ~~a national model run by ARPAE (<https://www.snpambiente.it/prodotti/previsioni-qualita-dellaria-in-italia/>).~~

The database of observed data used in this work, was built with the support of PREPAIR partners providing revised validated data for 2021.

3.5 CALPUFF simulation over Kosovo (FA5)

FA5 was operated by ATMOTERM Company between July 2020 and September 2022. Analyses were based on data available from Kosovo Air Quality Portal hosted by the Hydrometeorological Institute of Kosovo and Kosovo Open Data Platform (<https://airqualitykosova.rks-gov.net/en/>; <https://opendata.rks-gov.net/en/organization/khmi>). Forecast Service was using the following modelling tools: WRF meteorological prognostic model, CAMS ENSEMBLE Eulerian air quality models and CALPUFF Modelling System with 1 km receptor grid covering the Kosovo territory and 0.5 km grid applied in the main Kosovo cities. In addition a high resolution receptor network was created for Pristine, with the basic grid step of 200 m and 50 m along the roads. ~~Variable spatial grid size included almost 13000 grid points for Kosovo and about 12000 grid points for Pristine. Emission input data were prepared on the basis of Kosovo emission inventory for 2018 year (<https://ajri.niph.rks.org/>). Input data for WRF model were collected from the global forecasting model GFS NCEP with spatial resolution of 0.25°. Boundary condition for CALPUFF were obtained from Copernicus CAMS service (CAMS ENSEMBLE model).~~ The system includes an assimilation module implemented at the post-processing stage using available data from all monitoring stations in Kosovo. ~~FA5 presented 3-day forecast of four pollutants: PM₁₀, PM_{2.5}, NO₂ and O₃ and the European Air Quality Index (AQI). AQI calculation is based on EEA methodology (<https://airindex.eea.europa.eu/Map/AQI/Viewer/#>).~~

4 Results, Lessons learnt and Discussion

The proposed evaluation methodology for forecasting comes on top of the consolidated FAIRMODE protocol for assessment. The assessment MOO therefore comes first to provide a preliminary evaluation of the five forecasting applications (see Appendix C). This section ~~presents~~ focuses on the outcomes of applying the additional forecast objectives and criteria ~~validation of the five forecasting applications, focusing in particular~~ on the lessons learnt by ~~the~~ their application ~~of the proposed evaluation protocol~~ to different geographical contexts and spatial scales, ~~and~~ pointing out to the strengths and shortcomings of the approach ~~in highlighting the skills of forecasting applications.~~

4.1 MQO_f skills versus the capability of predicting Threshold Exceedances

~~The main assessment of the “fitness for purpose” is Forecast Modelling Quality Objective (MQO_f) outcomes are~~ presented here for three forecasting applications, covering different spatial scales, namely FA1 (European scale), FA2 (national scale), and FA4 (regional scale). Along with MQO_f outcomes, the skills of the three modelling applications in predicting threshold exceedances are provided as well. We present outcomes for PM10 daily mean and O₃ daily maximum of 8-hour average, since both indicators have a daily limit value set by the current European legislation (European Union, 2008).

~~Figs. 1-2, Figs. 23-4, and Figs. 3-5-6 show the outcomes for FA1, FA2, FA4 applications, respectively, in the form of panels of 4 plots. PM10 outcomes are on the left side of each panel provided in Figs. 1, 3, and 5 while Figs. 2, 4, and 6 present the O₃ ones are on the right side. MQI_f values are provided in the Forecast Target Plots (Janssen and Thunis, 2022), at the top of each panelFigure. Within these plots, MQI_f is represented by the distance between the origin and a given point (for each monitoring station). Values lower than 1 (i.e. within the green circle) indicate better capabilities than the persistence model (within the measurement uncertainty), whereas values larger than 1 indicate poorer performances. Indeed, the green area identifies the fulfilment of the MQO_f at each monitoring stations. The MQI_f associated to the 90th percentile worst station is reported in the upper left corner of the plots. This value is used as the main indicator in the proposed benchmarking procedure: its value should be less than or equal to 1 for the fulfilment of the benchmarking requirements. In other words, within the proposed protocol a forecasting application is considered “fit for purpose” if MQI_f is lower than 1 for at least 90% of the available stations. Note that passing the MQO_f test is intended here as a necessary condition for the use of the modelling results but it must not be understood as a sufficient condition that ensures that model results are of sufficient quality.~~

The outcomes of all categorical metrics included in the validation protocol are provided at the bottom of each panelFigure, by means of the Forecast Summary P-Normalized Reports. Within these plots, the statistical distribution (5th, 25th, 50th, 75th, 95th percentiles) of the outcomes of all the indicators defined in Sect. 2.2 are summarized and compared with the corresponding outcomes of the persistence model (i.e. the ratios of the skills are considered). Green area indicates that model performs better than the persistence model for that particular indicator.

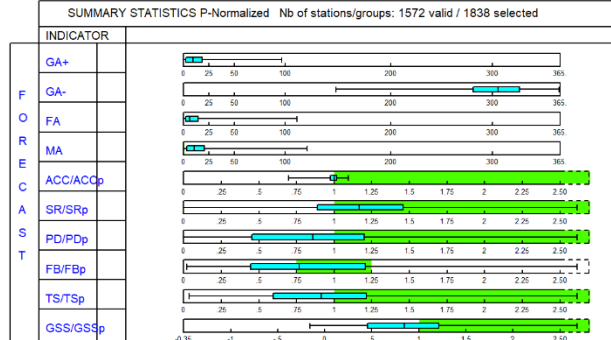
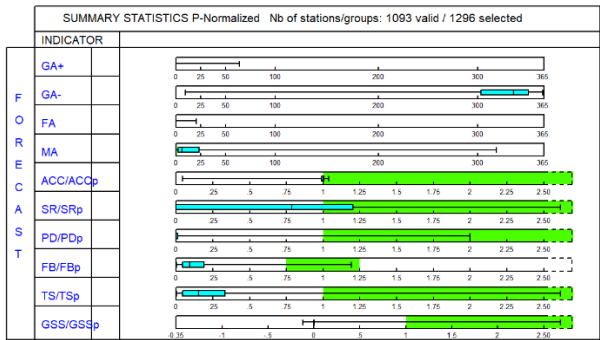
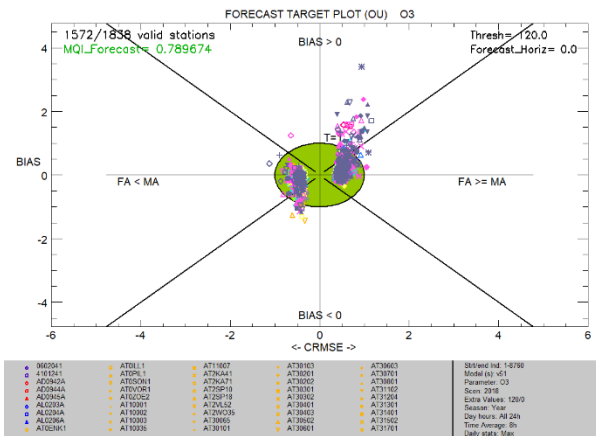
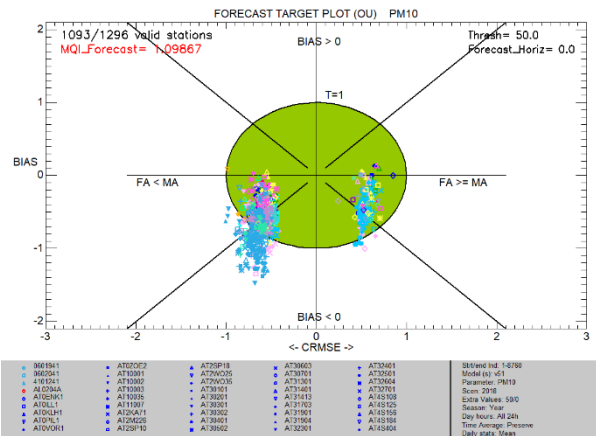


Figure 1: FA1 validation outcomes for PM10 (left) and O₃ (right). Forecast Target Plots (top) provide MQI_F values for each monitoring station, as the distance between the origin and a given point. Box plots in the Forecast Summary P-Normalized Reports (bottom) provide the statistical distribution (5th, 25th, 50th, 75th, 95th percentiles) of the categorical metrics.

330

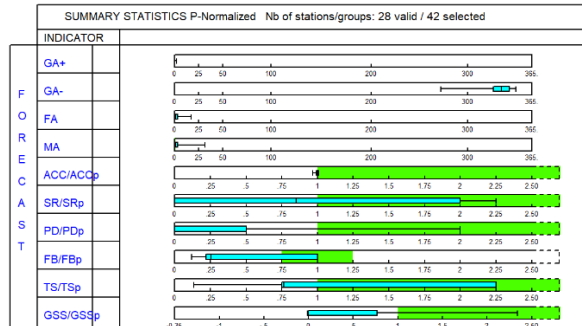
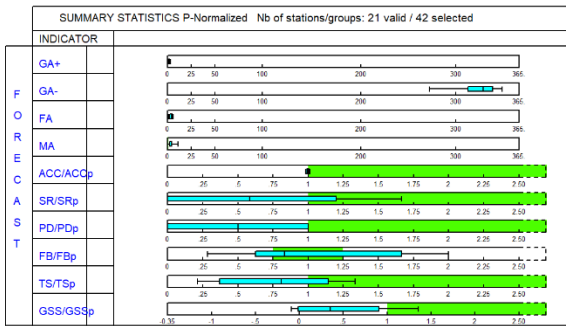
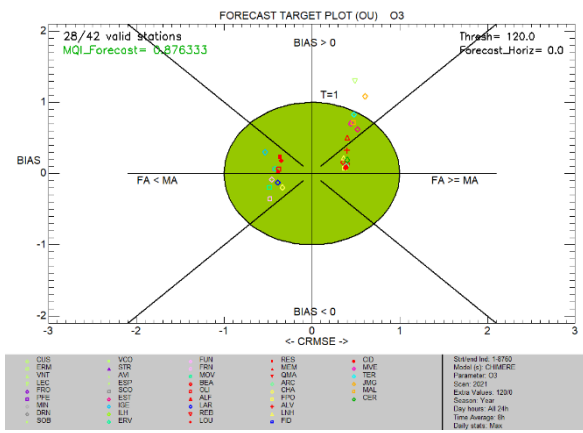
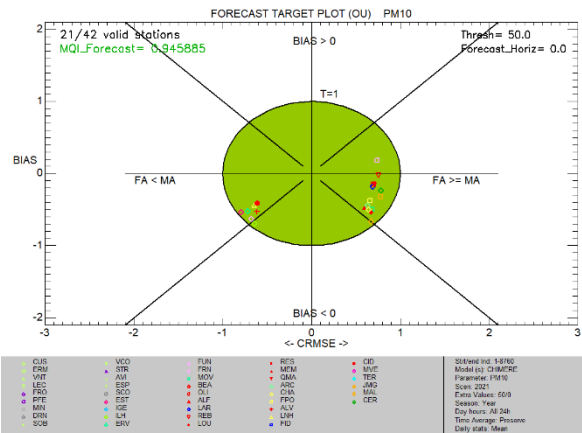


Figure 2: As Fig.1, for FA2 validation outcomes.

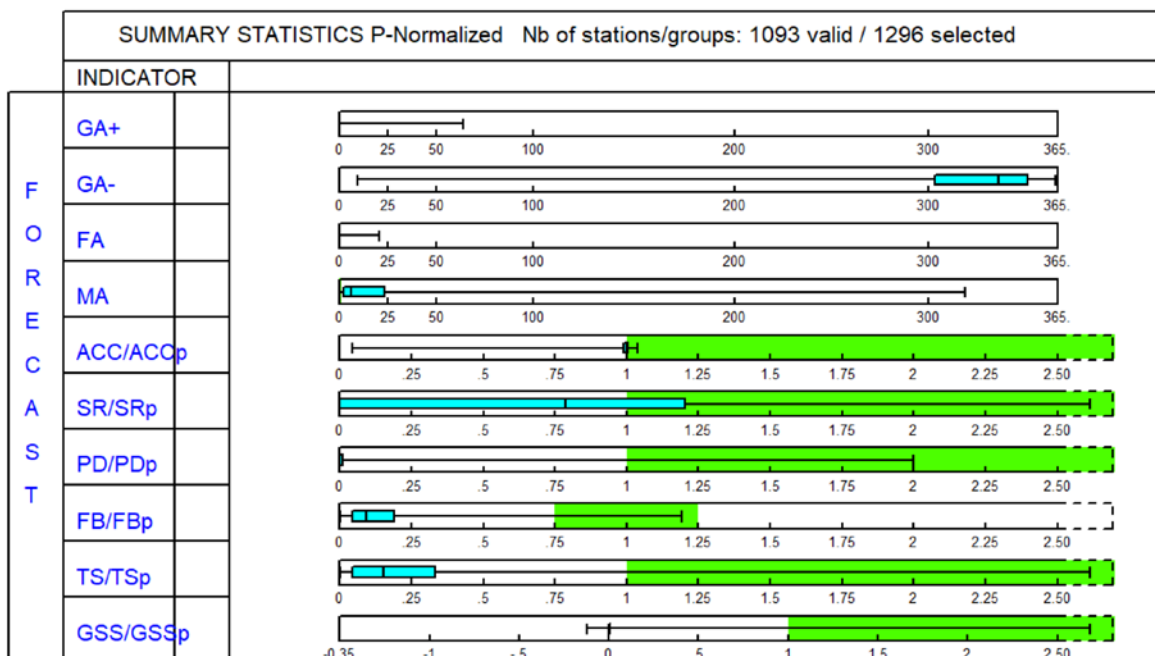
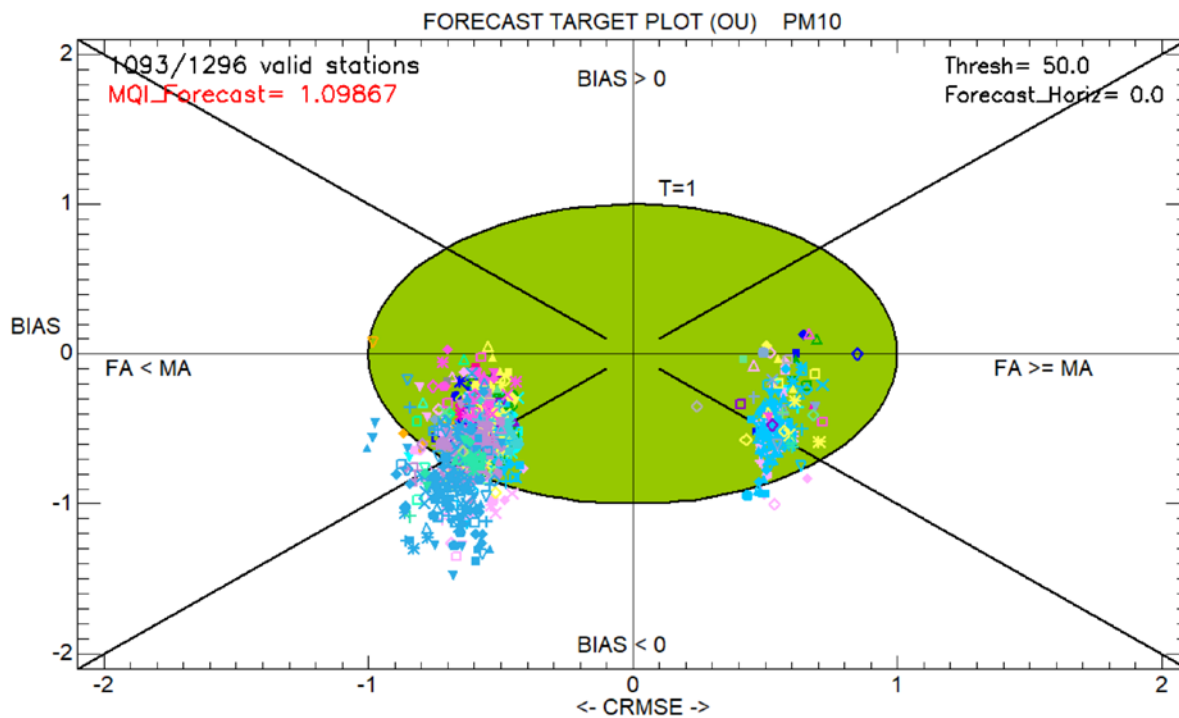
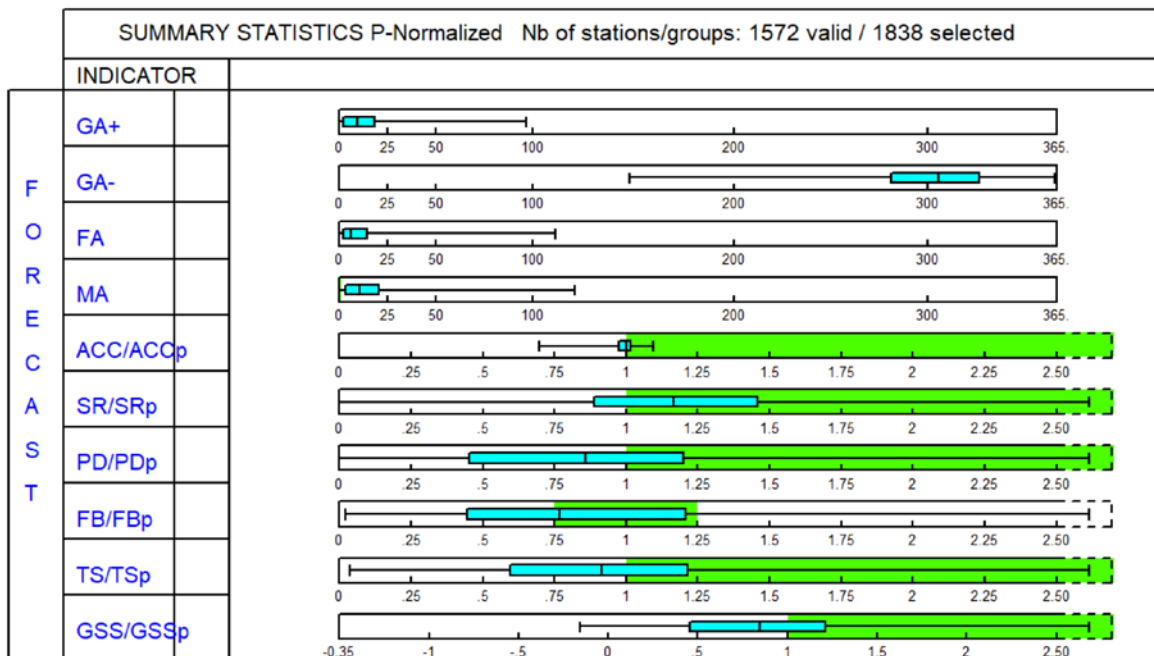
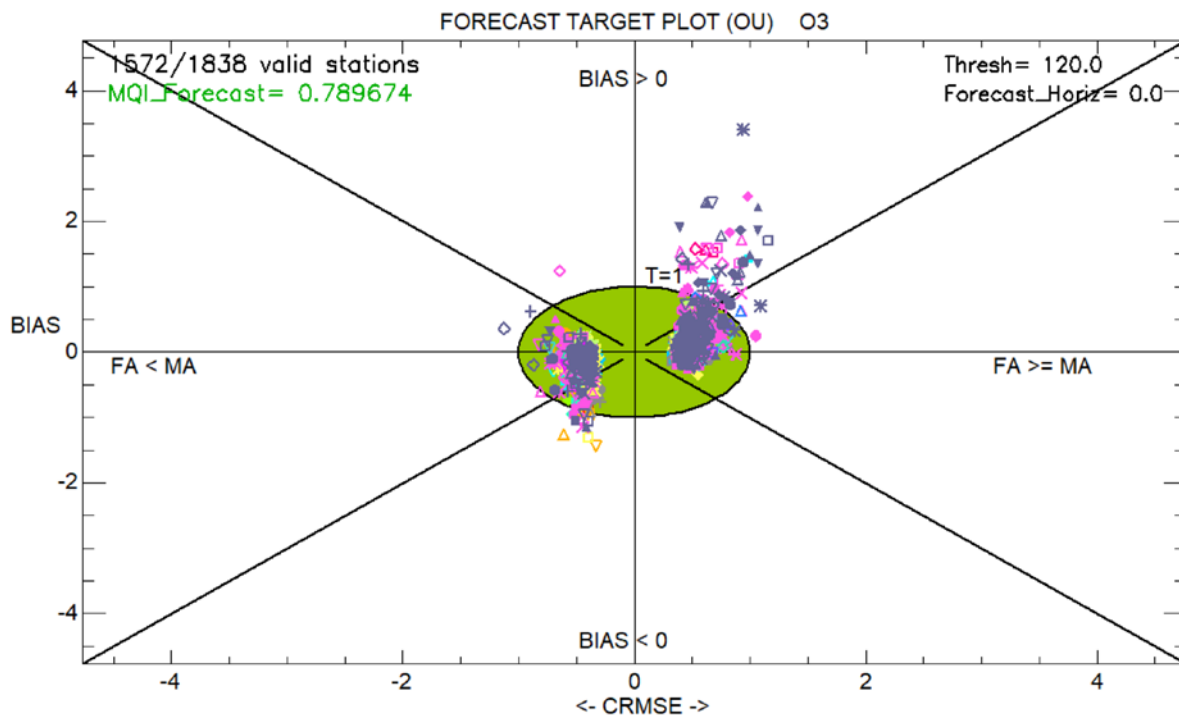


Figure 1: FA1 validation outcomes for PM10. Forecast Target Plots (top) provide MOI_f values for each monitoring station, as the distance between the origin and a given point. Box plots in the Forecast Summary P-Normalized Reports (bottom) provide the statistical distribution (5th, 25th, 50th, 75th, 95th percentiles) of the categorical metrics.



340 **Figure 2: FA1 validation outcomes for O₃. Forecast Target Plots (top) provide MOI_t values for each monitoring station, as the distance between the origin and a given point. Box plots in the Forecast Summary P-Normalized Reports (bottom) provide the statistical distribution (5th, 25th, 50th, 75th, 95th percentiles) of the categorical metrics.**

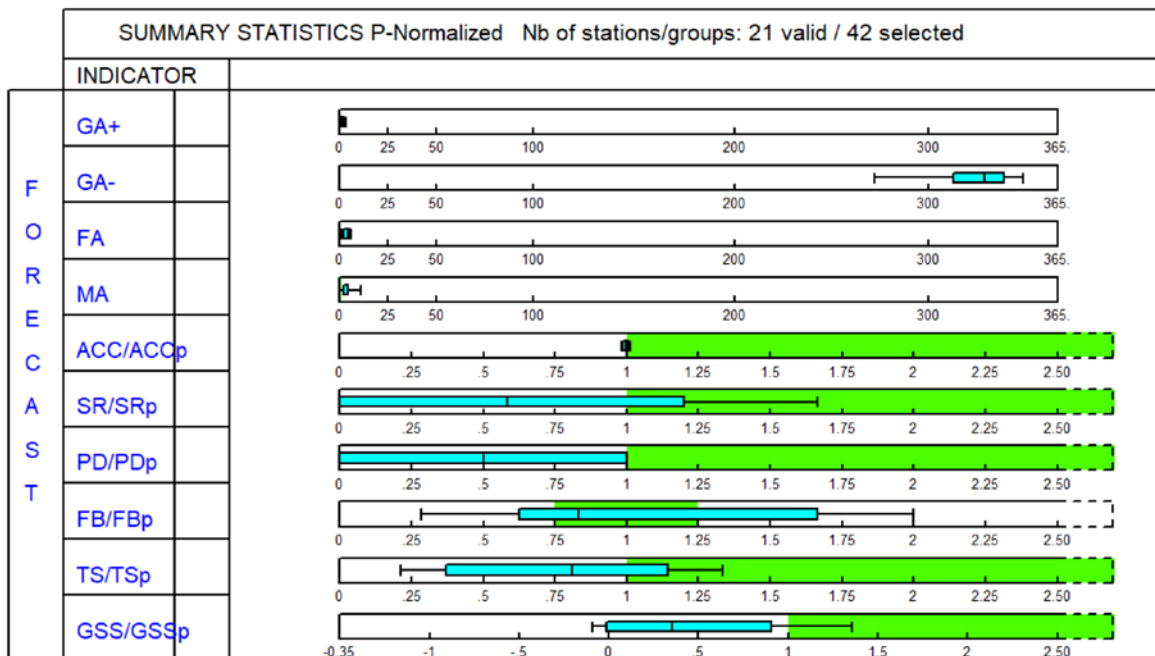
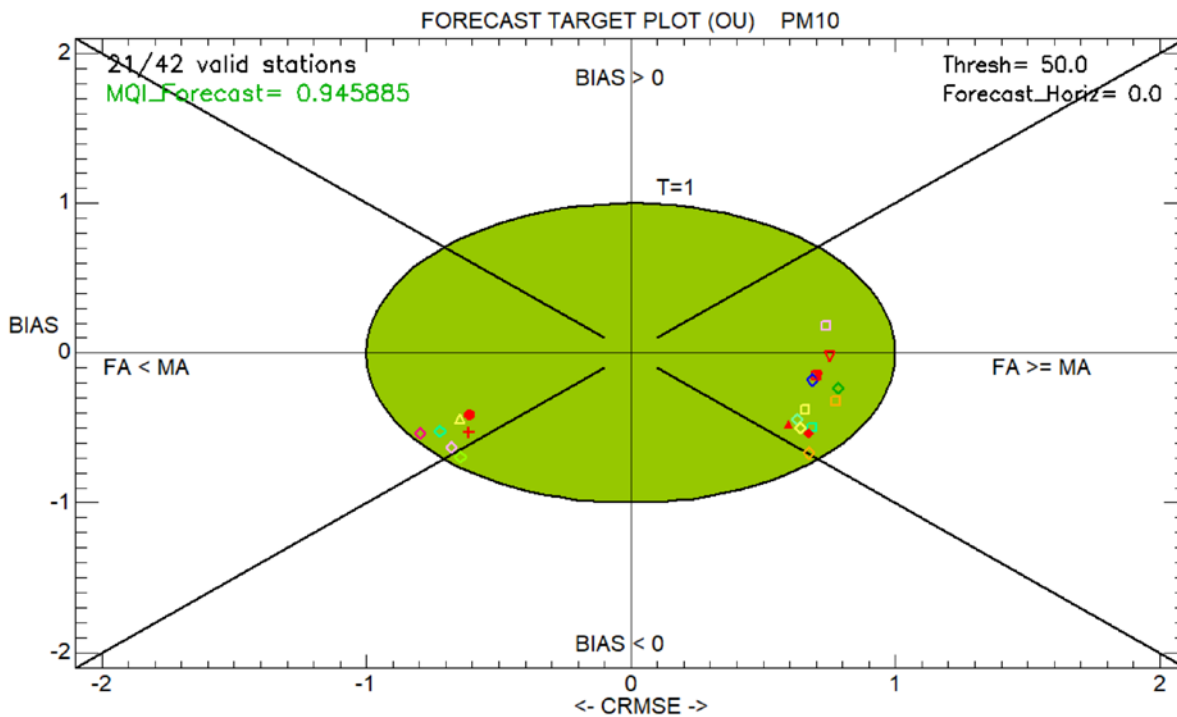


Figure 3: FA2 validation outcomes for PM10. Forecast Target Plots (top) provide MOI_f values for each monitoring station, as the distance between the origin and a given point. Box plots in the Forecast Summary P-Normalized Reports (bottom) provide the statistical distribution (5th, 25th, 50th, 75th, 95th percentiles) of the categorical metrics.

345

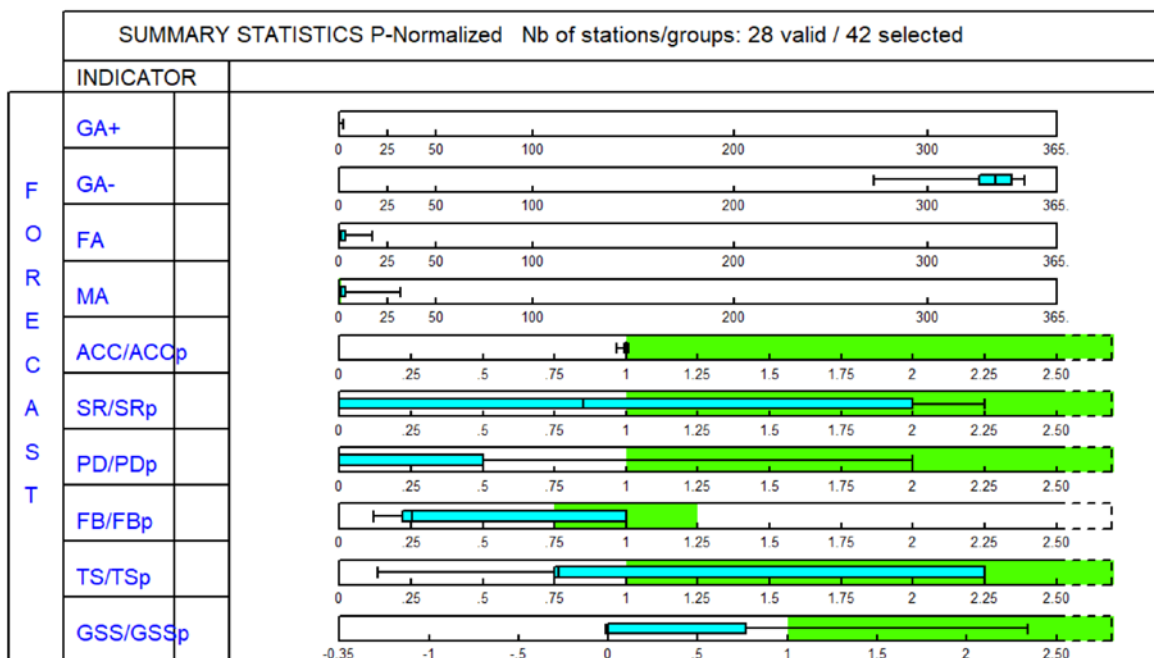
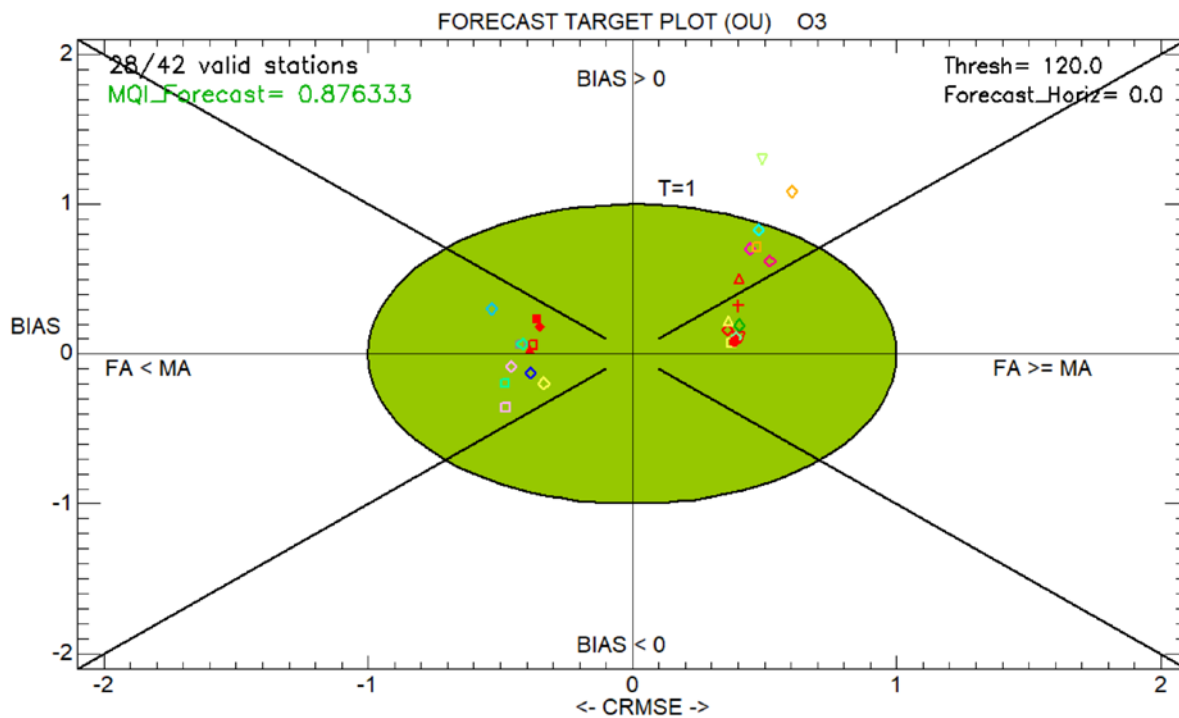


Figure 4: FA2 validation outcomes for O₃. Forecast Target Plots (top) provide MOI_F values for each monitoring station, as the distance between the origin and a given point. Box plots in the Forecast Summary P-Normalized Reports (bottom) provide the statistical distribution (5th, 25th, 50th, 75th, 95th percentiles) of the categorical metrics.

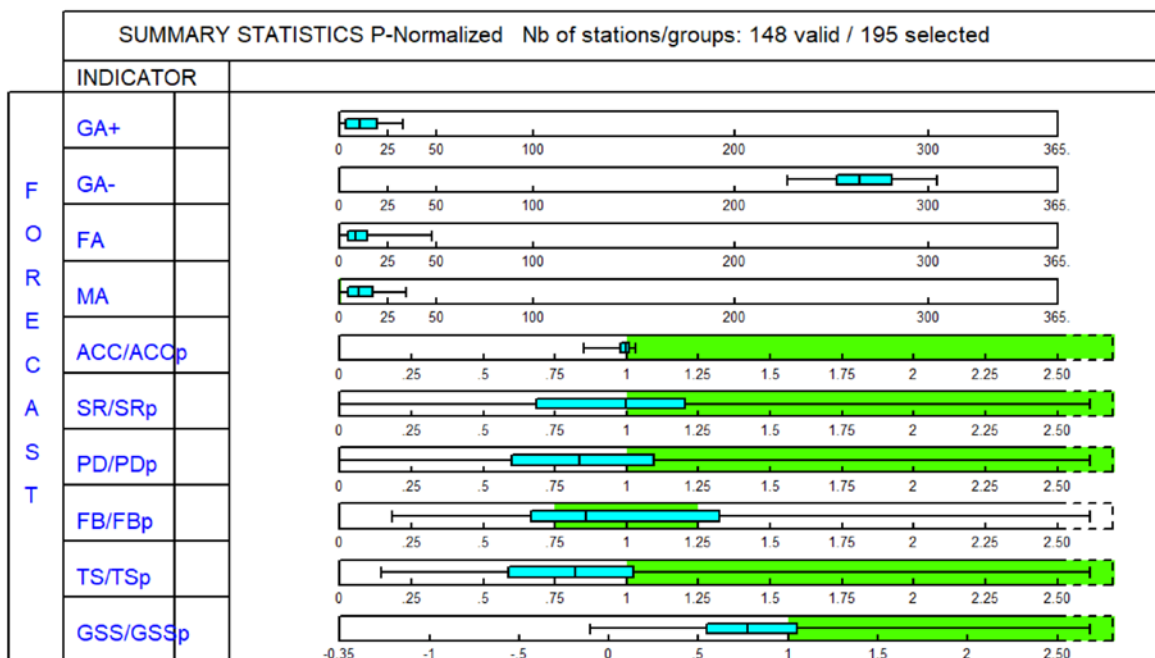
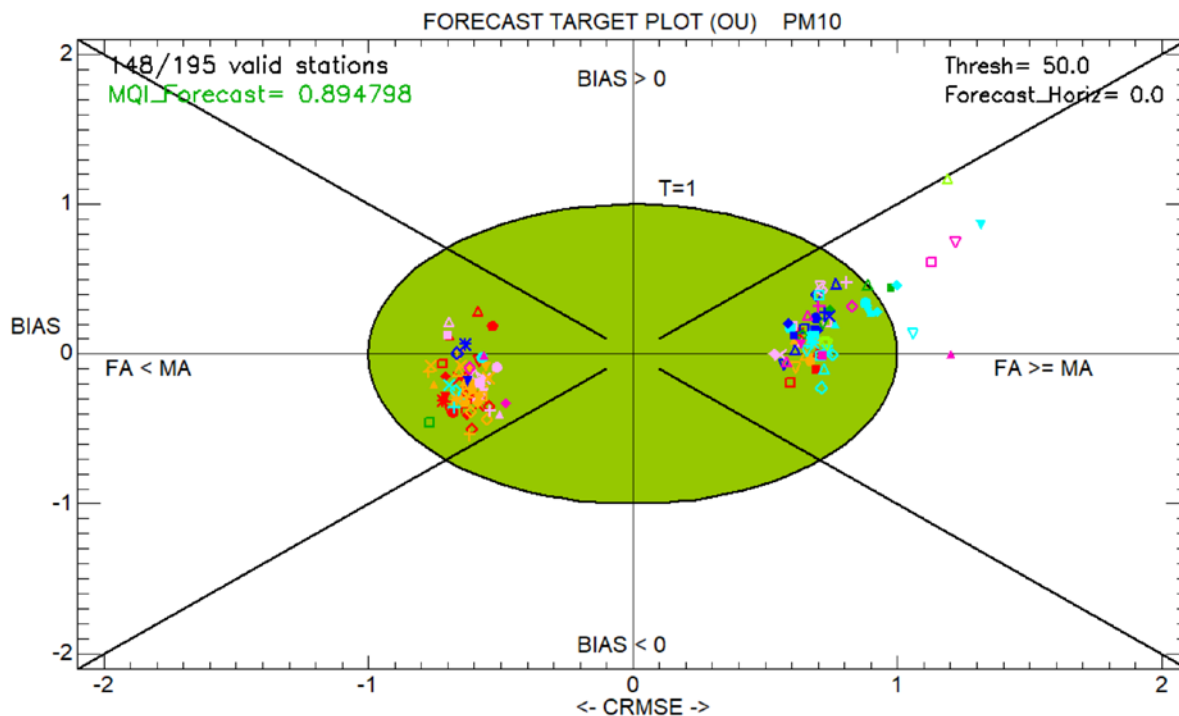
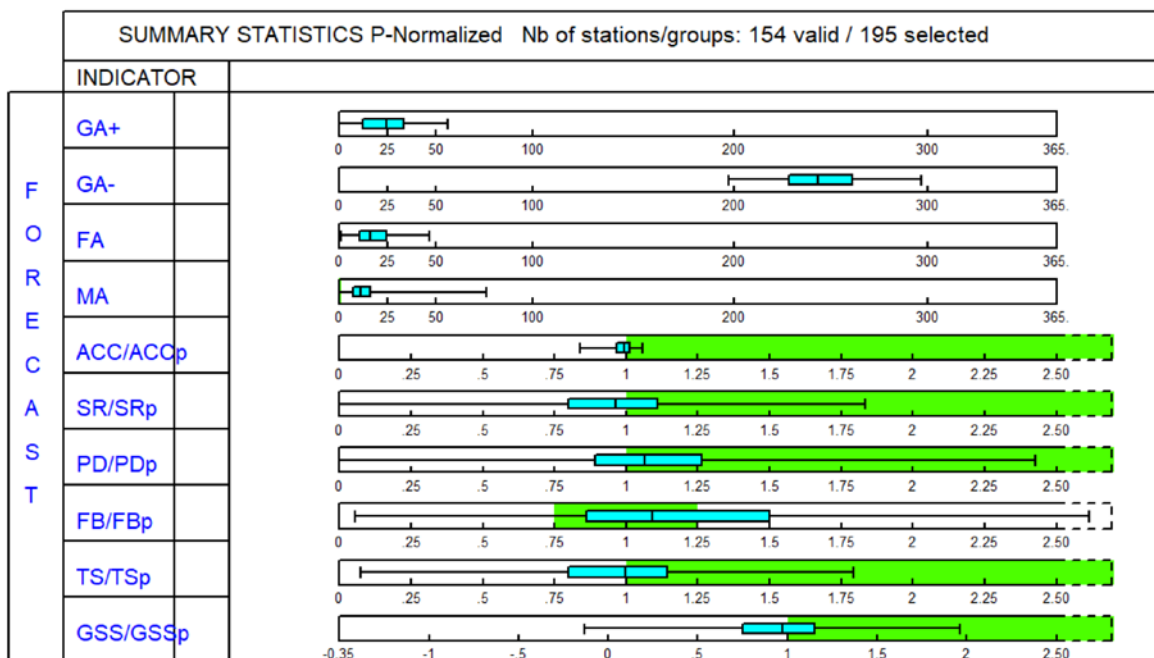
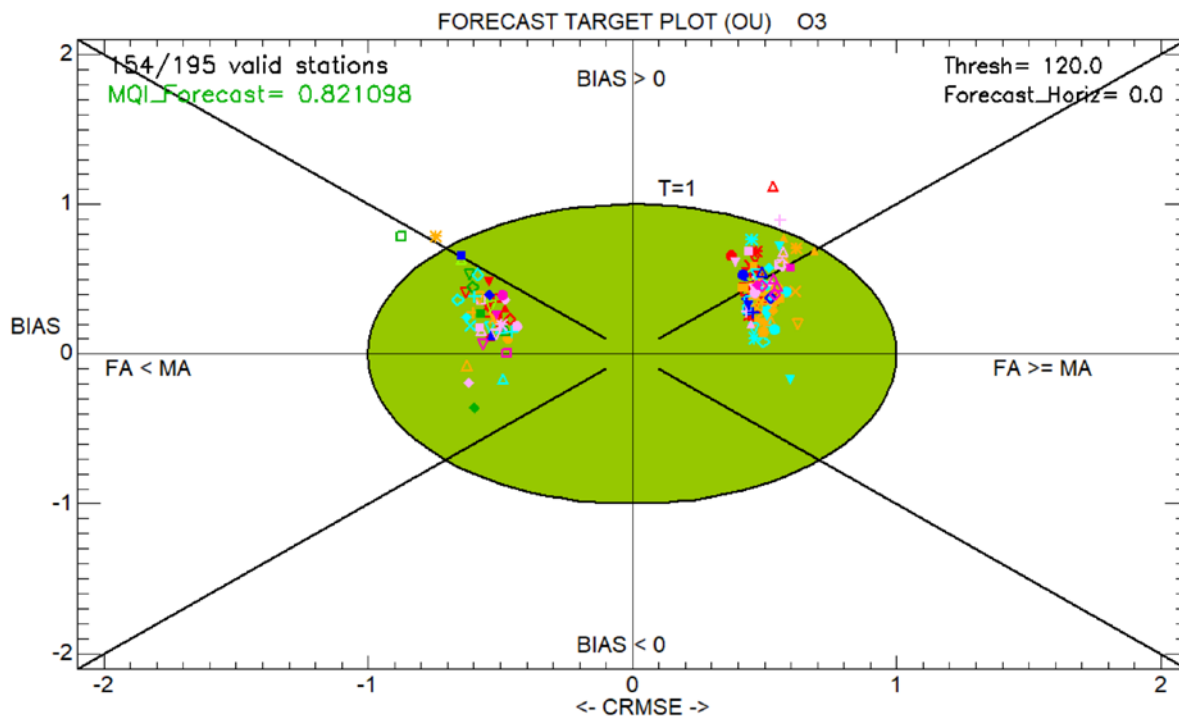


Figure 5: FA4 validation outcomes for PM10. Forecast Target Plots (top) provide MOI_f values for each monitoring station, as the distance between the origin and a given point. Box plots in the Forecast Summary P-Normalized Reports (bottom) provide the statistical distribution (5th, 25th, 50th, 75th, 95th percentiles) of the categorical metrics.



355

Figure 6: FA4 validation outcomes for O₃. Forecast Target Plots (top) provide MOI_f values for each monitoring station, as the distance between the origin and a given point. Box plots in the Forecast Summary P-Normalized Reports (bottom) provide the statistical distribution (5th, 25th, 50th, 75th, 95th percentiles) of the categorical metrics.

Forecast Target Plots outcomes indicate a very good level of quality of all forecast applications in simulating O₃. The 90th percentile of the MQI_f values is lower than 1 for all three forecast applications, indicating that model performs better than persistence in simulating O₃ at more than 90% of the available stations. FA2 and FA4 fulfil MQI_f requirements also in simulating PM10, instead there is room for improvement for the European scale simulation FA1 (90th percentile of the MQI_f values is slightly higher than 1). Further investigations show that most of the issues emerge in a limited part of the modelling domain (Turkey), where very high, and sometimes unlikely, PM10 values are measured at several monitoring sites for most of the year. Removing Turkish monitoring stations from the validation data set, MQO_f turns out to be fulfilled (Fig. ~~C1D1~~, in Appendix ~~CD~~). It is worth noting that the MQO_f outcomes are consistent with the standard assessment evaluation (Appendix C). Table C1 shows that the standard MQO is fulfilled for all O₃ forecast applications. For PM10, the MQI is higher than 1 but only for the FA1 simulation.

Concerning the capability in predicting exceedances, model performances improve moving from FA1 to FA4 applications (i.e. as spatial resolution increases) and skills are generally better in simulating O₃ than PM10. Concerning the comparison of the performances according to the different metrics, all forecast applications turn out to be better in avoiding false alarms than in reproducing all of them, since Success Ratio (SR) scores are generally better than Probability of Detection (PD) ones, especially for PM10.

In general, even if forecast applications are generally better than the persistence model according to the main outcome MQO_f (top plots of Figs. ~~1-6, Fig. 2, and Fig. 3~~), it becomes harder for them to beat the persistence model in predicting exceedances (bottom plots of Figs. ~~1-6, Fig. 2, and Fig. 3~~). Apart from few cases (namely the regional FA4 application), the median values of the statistical distribution of the outcomes are not in the green area, indicating that model performs worse than persistence at more than 50% of the available stations.

4.2 MPI Plot supporting the interpretation of MQO_f outcomes

When evaluating a forecasting application, it is of interest to assess the evolution of skills metrics with the forecast horizon. Indeed, a good forecasting application should not incur a substantial degradation of its performances along with forecast time.

FA3, carried out over Ireland by means of the OPAQ statistical system, was evaluated for each of the forecasted days, which included the current day (day0), tomorrow (day1) and the day after tomorrow (day2).

~~Fig. 4~~ In the following it is reported shows how performances in simulating PM10 vary along with the forecast days. More in detail outcomes for day0 (~~on the left~~) and day2 (~~on the right~~) are shown in Fig. 7 and Fig. 8, respectively. On the top of ~~the each panel~~ Figure, the Forecast Target Plots (described in the previous section) are reported. On the bottom, the Forecast MPI Plots are added, describing the fulfilment of both the criteria defined in Sect. 2.1 (i.e. MPI less than or equal to 1). Indeed, here the forecast performances (MFE_f) are compared to the persistence model performances (MFE_p) along Y axis (MPI_1) and to the Mean Fractional Uncertainty (MF_U) along the X axis (MPI_2). The green area identifies the area of fulfilment of both proposed criteria. The orange areas indicate where only one of them is fulfilled.

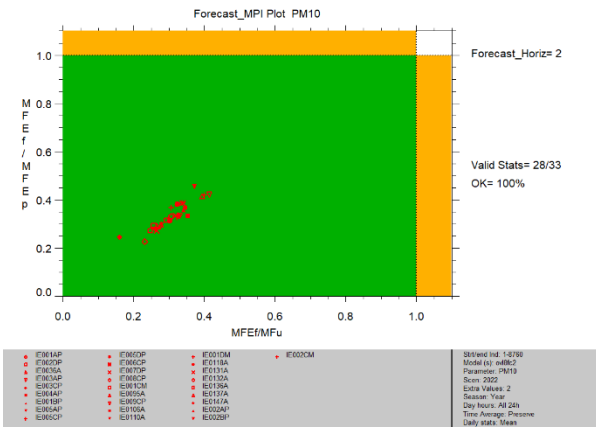
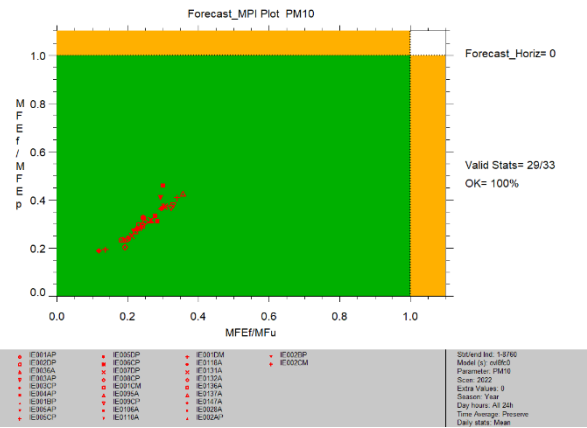
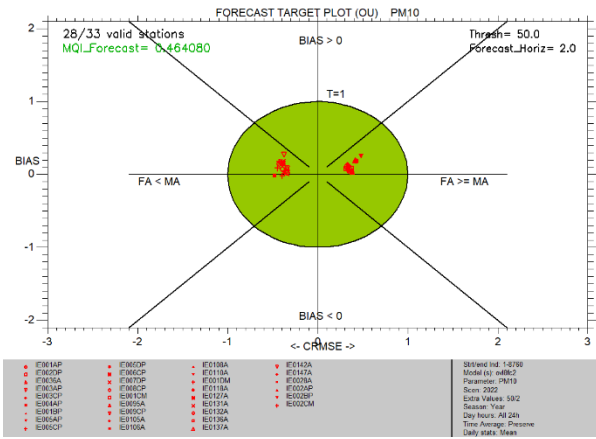
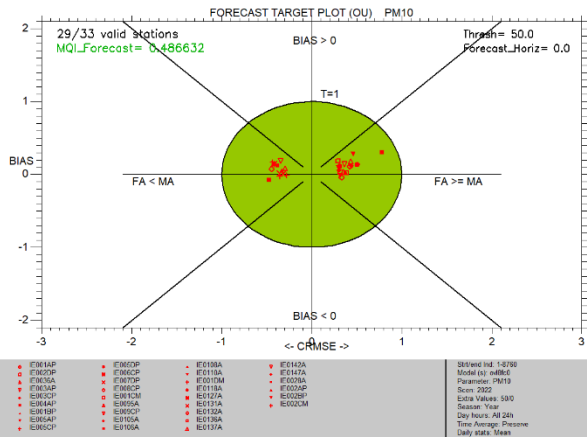


Figure 4: FA3 validation outcomes for day0 (left) and day2 (right). Forecast Target Plots (top) provide MQI_F values for each monitoring station, as the distance between the origin and a given point. Forecast MPI Plots (bottom) provide for each monitoring station MPI_1 along Y axis and MPI_2 along the X axis.

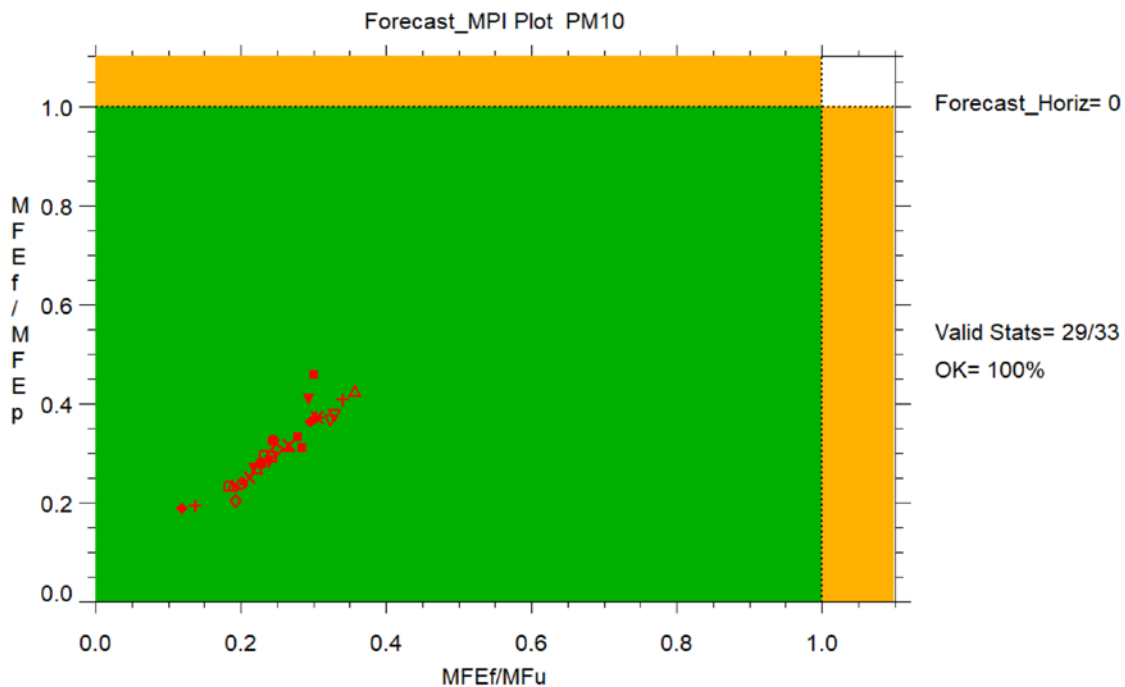
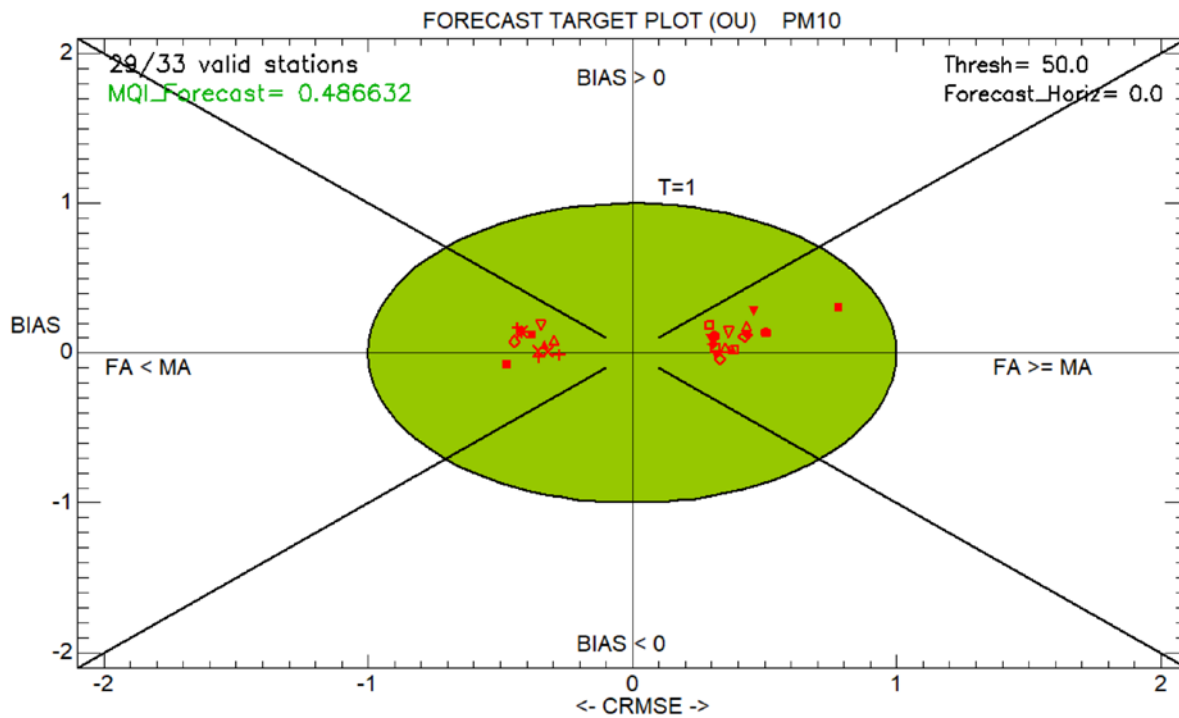


Figure 7: FA3 validation outcomes for day0. Forecast Target Plots (top) provide MOI_f values for each monitoring station, as the distance between the origin and a given point. Forecast MPI Plots (bottom) provide for each monitoring station MPI_f along Y axis and MPI_p along the X axis.

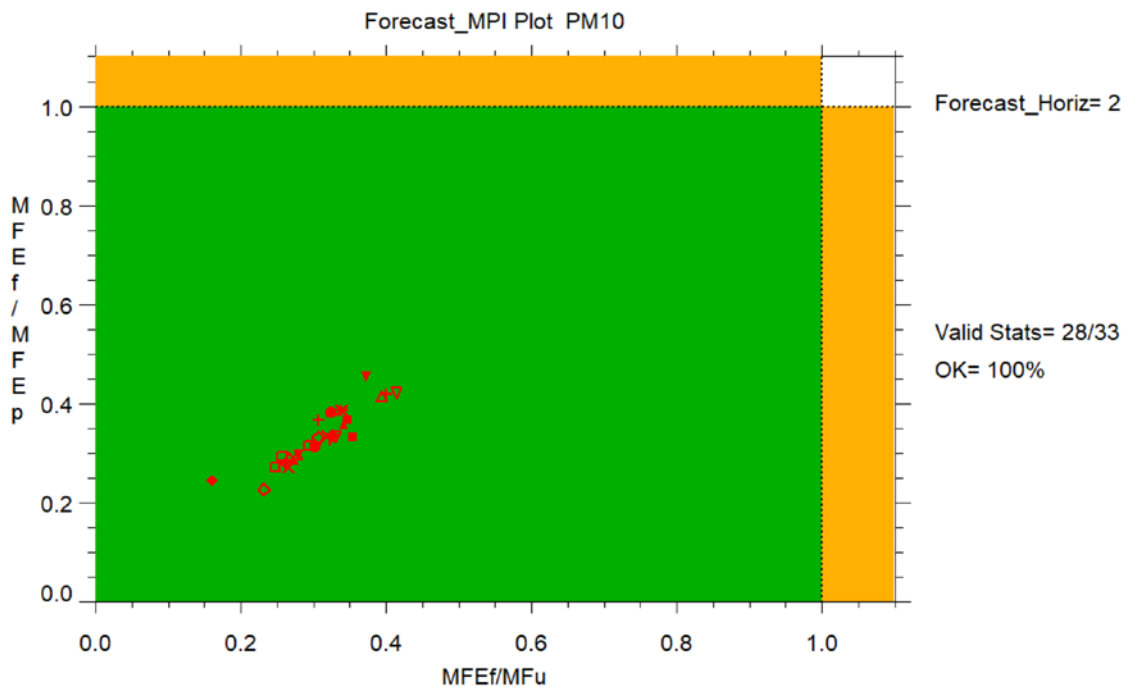
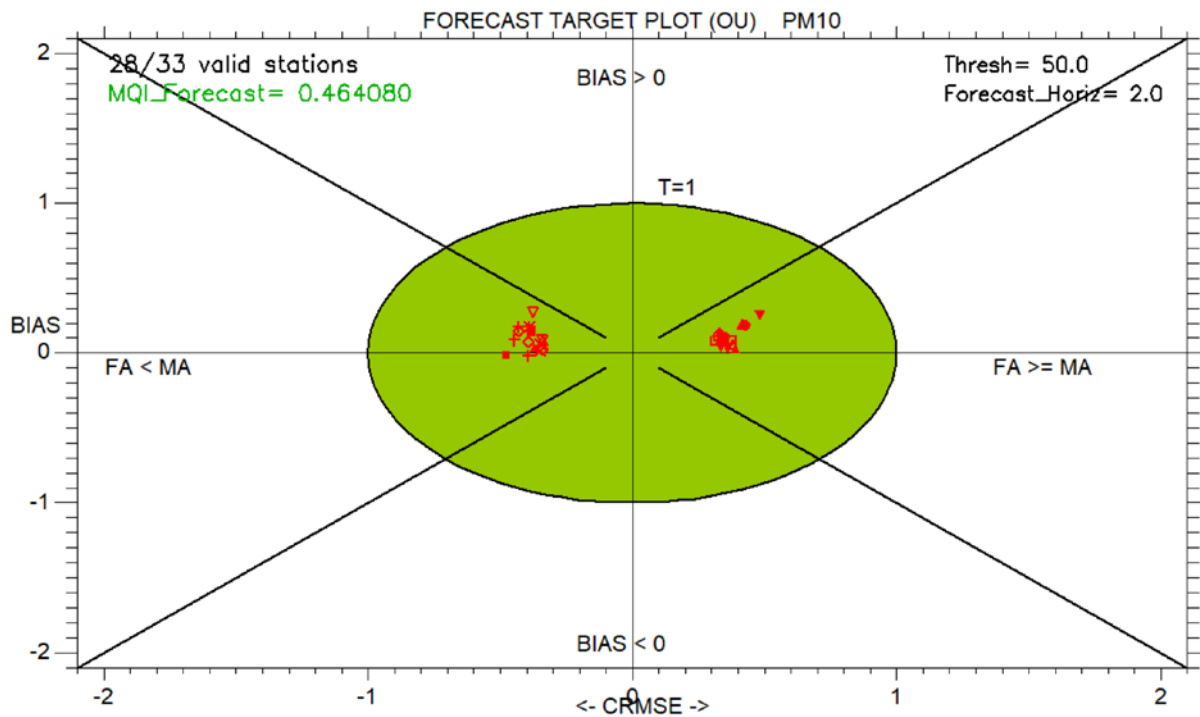


Figure 8: FA3 validation outcomes for day2. Forecast Target Plots (top) provide MOI_f values for each monitoring station, as the distance between the origin and a given point. Forecast MPI Plots (bottom) provide for each monitoring station MPI_f along Y axis and MPI_p along the X axis.

405 Outcomes in Figs. 4-7-8 indicate a very good level of quality of the forecast application, since Modelling Quality Objective is fulfilled (top), together with the two additional Performance Criteria (bottom). These outcomes are consistent with the standard *MQO* skills provided in Table C1 of Appendix C, which points out very good performances of FA3 for PM10, namely the best performances among all forecast applications.

Concerning the evolution of skills metrics with forecast horizon, according to the Forecast Target Plot outcomes (top), 410 modelling performances unexpectedly get better from day0 to day2, since the MQI_f value associated to the 90th percentile worst station (reported in the upper left corner of the plots) turns out to get lower. According to the Forecast MPI Plots (bottom), performances remain almost constant with forecast horizon, indicative of a good behaviour of the modelling application. Moreover, Forecast MPI Plots help to clarify that the unrealistic improvement of model performances from day0 to day2, pointed out by the Forecast Target Plots, is due to persistence model performances degradation. Indeed, moving 415 from day0 to day2, the forecast model performances get slightly better along Y axis, where they are normalized to persistence model skills, but they slightly deteriorate along X axis, where they are considered regardless of persistence aspects. In other words, model performances slightly deteriorate along with the forecast days but persistence model deteriorate more, so that performances ratios (i.e. both MQI_f and MPI_l values) get lower.

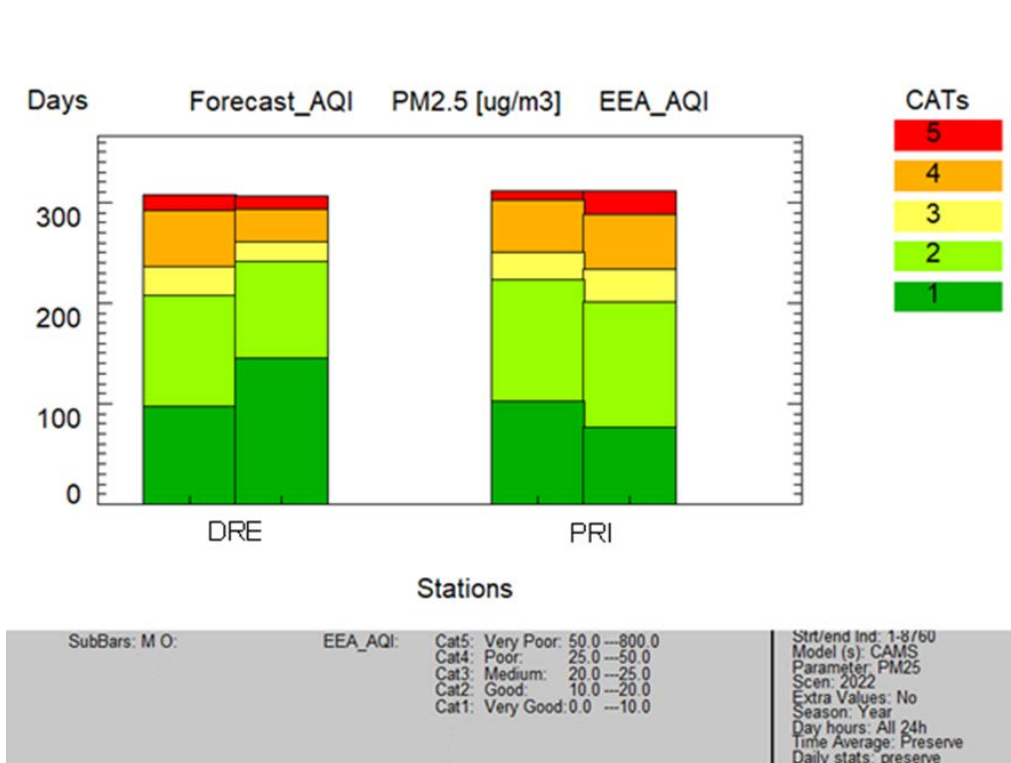
4.3 Assessment of modelling application capability in predicting Air Quality Indices

420 The current approach for the assessment of modelling application capability in predicting Air Quality Indices is based on a cumulative analysis for answering the following questions: “Are citizens correctly warned against high pollution episodes?” or in another words: “Does the model properly forecast AQI levels?”

Air Quality Indices are designed to provide information on local air quality. Moreover, within the proposed validation protocol, the capability of correctly predicting AQI is assessed at single monitoring stations. For these reasons, FA5 at the 425 local scale is the most suitable for testing the proposed approach. Indeed, it was carried out at high spatial resolution and focuses on only two monitoring sites, located in two cities in Kosovo: Pristine (the capital) and Drenas.

Before analysing AQI results for PM2.5, it has to be mentioned that the FA5 standard *MQO* is fulfilled for all available pollutants (Table C1 in Appendix C). Concerning additional features of the forecasting validation protocol, both the Forecast Target Plot and the Forecast MPI Plot show very good performances for both locations. The Forecast Summary P- 430 Normalized Report indicates good model performance in Drenas and some room for improvements in Pristine location due to underestimation of PM2.5 episodes.

Fig. 5-9 provides the AQI diagram, based on EEA classification, for PM2.5 and for day0 forecast. For each station, the bar plot shows two paired bars: the number of predicted (left bar) and measured (right bar) concentration values that fall within a given air quality category. In Drenas, forecast values populate categories 2 (“Good”), 3 (“Medium”), and 4 (“Poor”) to a 435 greater extent than the measurements. On the contrary, in Pristine forecast values are more frequent than the measurements at the lowest AQI (“Very Good”).



440 **Figure 59:** FA5 validation outcomes for PM2.5 at Drenas and Pristina. AQI diagram provide for each monitoring station the number of predicted (left bar) and measured (right bar) concentration values that fall within each air quality category. The last two EEA AQI classes (“Very poor” and “Extremely poor”) are merged into one.

Overall, Fig. 5-9 points out that FA5 generally overestimates PM2.5 concentration levels in Drenas and underestimates them in Pristine. Anyway, AQI forecast bar plots provide information about the total number of occurrences in each AQI class but there is no information about the correct timing of the forecasted AQI level.

445 So, there is room for future improvement and other additional outputs could be included within the protocol. In particular, multi-category contingency tables can be created for each station and multi-categorical skill scores can be computed, according to literature (e.g. EPA, 2003). Outcomes can be plotted for single stations or describing, for each AQI class, skill scores statistical distribution among the stations.

450 For example, in Fig. 6-10 an in-depth insight of AQI assessment is proposed for Drenas (top) and Pristina (bottom). Two additional multi-categorical metrics are proposed. Both of them are computed for each AQI level and are based on the comparison between forecast and measurement values considering also the correct timing of the predicted AQI level. *AQI comparability* (left plots in Fig. 6-10) represents, for each of the five AQI classes, the percentage of the correct forecast events in this class with respect to the total events based on measurements. Since *AQI comparability* values are percentages, they range from 0 to 100, being 100 the optimal value. *TS_AQI* (right plots in Fig. 6-10) is computed according to the same

definition of *TS* in Table 1. Indeed, here multiple thresholds (i.e. class limits) are taken into account and so multiple
 455 outcomes, one for each AQI class, are provided. *TS_AQI* values range from 0 to 1, being 1 the optimal value.

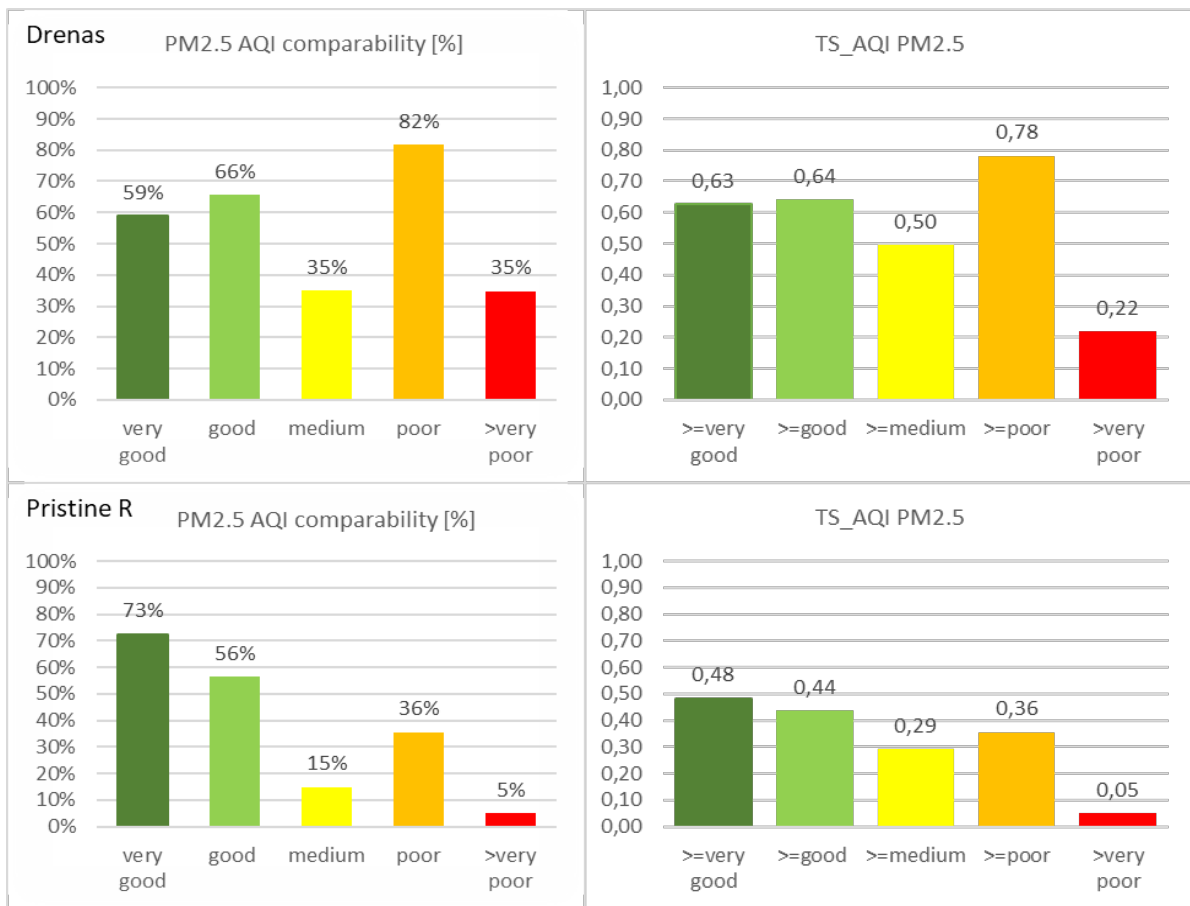


Figure 610: Multi-categorical metrics outcomes for Drenas (top) and Pristina (bottom). AQI comparability plots (left) provide for each AQI class the percentage of the correct forecast events with respect to the total events based on measurements. *TS_AQI* plots (right) provide for each AQI class *TS_AQI* values.

460 *AQI Comparability* and *TS_AQI* in Fig. 6-10 provide additional information with respect to AQI diagram. For example, in
 the case of Drenas, it turns out that, according to both the metrics, the best agreement between forecast and measurements in
 predicting the correct timing of the occurrences are found for “Poor” AQI class. It is also worth noting that, even if
 according to cumulative analysis (Fig. 59) forecast and measurements present a similar number of occurrences in both the
 “Medium” and the “Very Poor” classes, according to *AQI Comparability*, these classes are characterized by the worst
 465 performances. *TS_AQI* gives additional information about the model performances, which is especially noticeable for the
 “Medium” and “Very poor” classes, as it defines levels differently (“Medium” class means “Medium” class and all higher
 classes – “Poor” and “Very poor”). In this case the “Medium” class is characterized by better performances than the “Very

Poor” class. In the case of Pristine location, the best performances, according to both the metrics, are achieved for low concentrations (“Very Good” and “Good” classes) and the worst ones for “Very Poor” and “Medium” AQI levels. It is also worth noting that the best agreement is found for “Good” class, according to the cumulative comparison (Fig. 59), but it is better for “Very Good” class if the timing of the occurrences is taken into account (Fig. 610).

4.4 Discussion

Several lessons were learnt from the results presented here. The main ~~“fitness for purpose”-proposed~~ criterion (MQO_f) turned out to be useful for ~~a comprehensive evaluation of~~ evaluating the strengths and the shortcomings of a forecasting ~~evaluation~~ application, focusing on features which could not be addressed with the assessment evaluation approach.

Side outcomes, included within the protocol, can help in deepening the analysis. For example, MPI s analysis based on MFE helps in interpreting the outcomes, since MPI_2 is formulated regardless of persistence aspects, providing ~~as an added value,~~ details an evaluation of on the model performances ~~quality itself.~~

Consistently with the FAIRMODE approach, the measurement uncertainty is considered within the MQO_f formulation. While values are currently based on maximum uncertainties (95th percentile), these could be modified in the future to obtain a consensus level of stringency for the MQO_f , i.e. a level reachable for best applications while stringent enough to preserve sufficient quality. In Appendix E the outcomes of a sensitivity analysis are provided, in which we investigate the impact of the value chosen as representative measurement uncertainty.

Concerning the capability in predicting the exceedances, it turned out that, regardless of the spatial scale and the pollutants, even if a forecast application is better than the persistence model according to the general ~~“fitness for purpose”-evaluation~~ criterion (MQO_f), it can be worse in correctly providing categorical answers. Indeed, the difficulty in beating the persistence model skills is not infrequent in weather forecasting applications (Mittermaier, 2008). Moreover, it is worth noting that, differently from MQO_f analysis, the evaluation of the model capability in predicting the exceedances, being based on the definition of fixed thresholds, does not take into account the measurement uncertainty. For these reasons, a “fitness for purpose” criterion concerning exceedances metrics (e.g., which percentiles of a categorical indicator should be in the green area in order to define its skill “good enough”? and following on from that, how many indicators should be “good enough” in order to define the forecast application “fit for purpose”?) is not definitively set within the proposed protocol. Indeed, some more discussion based on further tests on forecasting applications are needed.

The greatest room for improvement turns out concerning the evaluation of the capability of the forecasting application in predicting AQI levels. The current approach is based on a cumulative analysis and no information is provided about the correct timing of the forecasted AQI levels. To account for this, some preliminary tests were carried out based on two additional multi-categorical metrics, which sound interesting in complementing the current approach. The main weakness of the proposed approach is the large number of different values to be provided, so making this type of outcome usable only for single monitoring stations. Moreover, the question of which level of performances in AQI predicting is “good enough” is currently an open issue and benchmarking of several forecasting applications is needed to establish some quality criteria.

5 Conclusions

A standardized validation protocol for air quality forecast applications was ~~proposed~~made available, following FAIRMODE community discussions on how to address specific issues typical of forecasting applications.

505 A-The proposal of a common benchmarking framework for model developers and users supporting policymaking under the European Air Quality Directives is a major achievement.

The proposed validation protocol enables an objective assessment of the “fitness for purpose” of a forecasting application, since it relies on the usage of a reference forecast as a benchmark (i.e. the persistence model), includes the measurement uncertainty, and bases the evaluation on the fulfilment of specific performance criteria, defining an acceptable quality level of the given model application. On top of a pass/fail test to ensure fitness-for-purpose (intended as a necessary but not
510 sufficient condition), a series of indicators is proposed to further analyse the strengths and weaknesses of the forecast application.

Moreover, relying on a common standardized validation protocol, the comparison of performances of different forecast applications, within a common benchmarking framework, is made available.

The application of the methodology to validate several forecasting simulations across Europe, using different modelling
515 systems and covering various geographical contexts and spatial scales, suggested some general considerations about its usefulness.

The main “fitness for purpose” criterion, describing the global performances of the model application with respect to persistence skills, proves to be useful for a comprehensive evaluation of the strengths and the shortcomings of a forecasting application. Generally, the forecast Modelling Quality Objective turns out to be achievable for most of the examined
520 validation exercises. When the criterion was not addressed, side analyses and outcomes, included within the protocol, helped in deepening the analysis and in identifying the most critical issues of the forecasting application.

On the other hand, it turned out that, regardless of the spatial scale and the pollutants, it can be hard for a forecast application beating persistence model skills in correctly providing categorical answers, namely on exceedances of concentration thresholds. Therefore, further tests and analyses are needed in order to provide some criteria for defining the “fitness for
525 purpose” of a forecasting application in predicting exceedances.

The last model capability assessed within the proposed validation protocol concerns the correct prediction of Air Quality Indices, designed to provide citizens with effective and simple information about air quality and its impact on their health. The current approach is based on a cumulative analysis of relative distributions of observed and forecasted AQIs. As no information is provided about the correct timing of the forecasted AQI levels, further developments are foreseen based on
530 multi-category contingency tables and multi-categorical skill scores.

Actually, discussion on the proposed approach will go on within the FAIRMODE community and upgrades and improvements of the current validation protocol will be fostered by its usage. In particular, it will be of interest to collect feedback from in-depth diagnostic analyses focusing on the validation of specific forecast applications, using both the

535 proposed criteria and the threshold-based categorical metrics to gain further insights. Anyway, from its preliminary applications across Europe, the methodology turns out to be sufficiently robust. Indeed, both the methodology and the software are publicly available for testing and application, especially targeting air quality forecasting services supporting policymaking in European Member States.

Appendix A

Measurement uncertainty $U(O_i)$ as a function of the concentration values O_i can be expressed as follow:

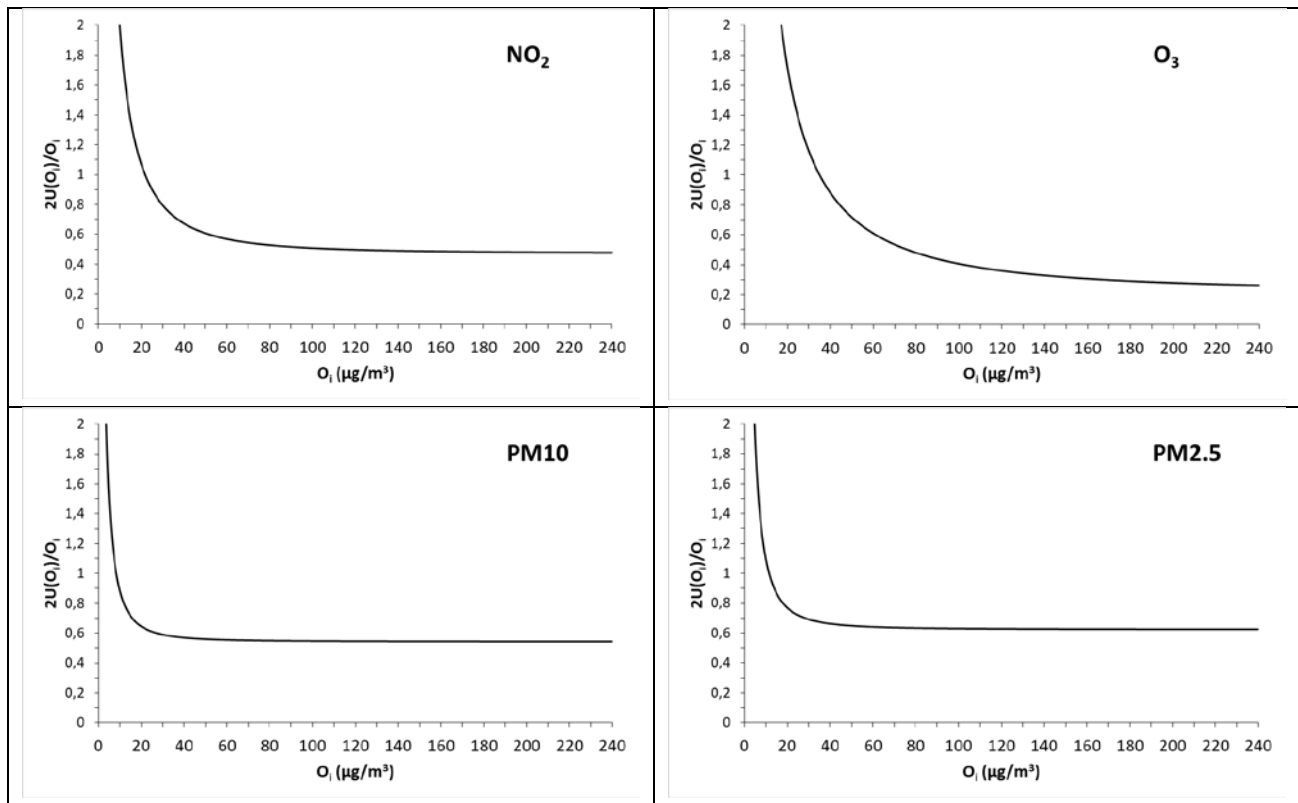
$$540 \quad U(O_i) = U_r(RV) \sqrt{(1 - \alpha^2)O_i^2 + \alpha^2 RV^2} \quad (A1)$$

545 An in-depth description of the rationale and the formulation of the measurement uncertainty estimation is provided in Thunis et al. (2013) and Pernigotti et al. (2013) for O₃, and PM and NO₂, respectively. More in details, the formulation of the measurement uncertainty as a function of the measured concentration is based on two coefficients: $U_r(RV)$, i.e. the relative uncertainty around a reference value RV and α , i.e. the fraction of uncertainty non proportional to the concentration value. It is important to note that we use as representative for the measurement uncertainty the 95th percentile highest value among all uncertainty values. For PM10 and PM2.5 the results of a JRC instrument inter-comparison (Pernigotti et al., 2013) have been used whereas a set of EU AIRBASE stations available for a series of meteorological years has been used for NO₂ and analytical relationships have been used for O₃. These 95th percentile uncertainties only include the instrumental error. Parameters $U_r(RV)$, RV , and α for $U(O_i)$ calculation for NO₂, O₃ and PM are provided in Table A1.

550

Table A1 Parameters for the calculation of measurement uncertainty.

	$U_r(RV)$	RV	α
NO₂	0.24	200 µg/m ³	0.20
O₃	0.18	120 µg/m ³	0.79
PM10	0.28	50 µg/m ³	0.25
PM2.5	0.36	25 µg/m ³	0.50



555 **Figure A1:** Double relative measurement uncertainties as a function of concentration values for NO_2 (top, left), O_3 (top, right), PM_{10} (bottom, left) and $\text{PM}_{2.5}$ (bottom, right).

Appendix B

Table B1 Contingency Table.

Forecast Events	Yes	FA	GA+
	No	GA-	MA
CONTINGENCY TABLE		No	Yes
		Observed Events	

Appendix C

560 The standard Modelling Quality Objective (MOQ), valid for assessment, is defined by the comparison of model-observation differences (namely, the root mean square error, RMSE) with a quantity proportional to the measurement uncertainty.

$$MQI = \frac{RMSE}{\beta \sqrt{\frac{\sum_{i=1}^N (U(O_i))^2}{N}}} \quad (C1)$$

β is set to 2, thus allowing the deviation between modelled and observed concentrations to be twice the measurement uncertainty. Measurement uncertainty $U(O_i)$ as a function of the concentration values O_i is defined in Appendix A.

565 The MOQ is fulfilled when MQI is less than or equal to 1.

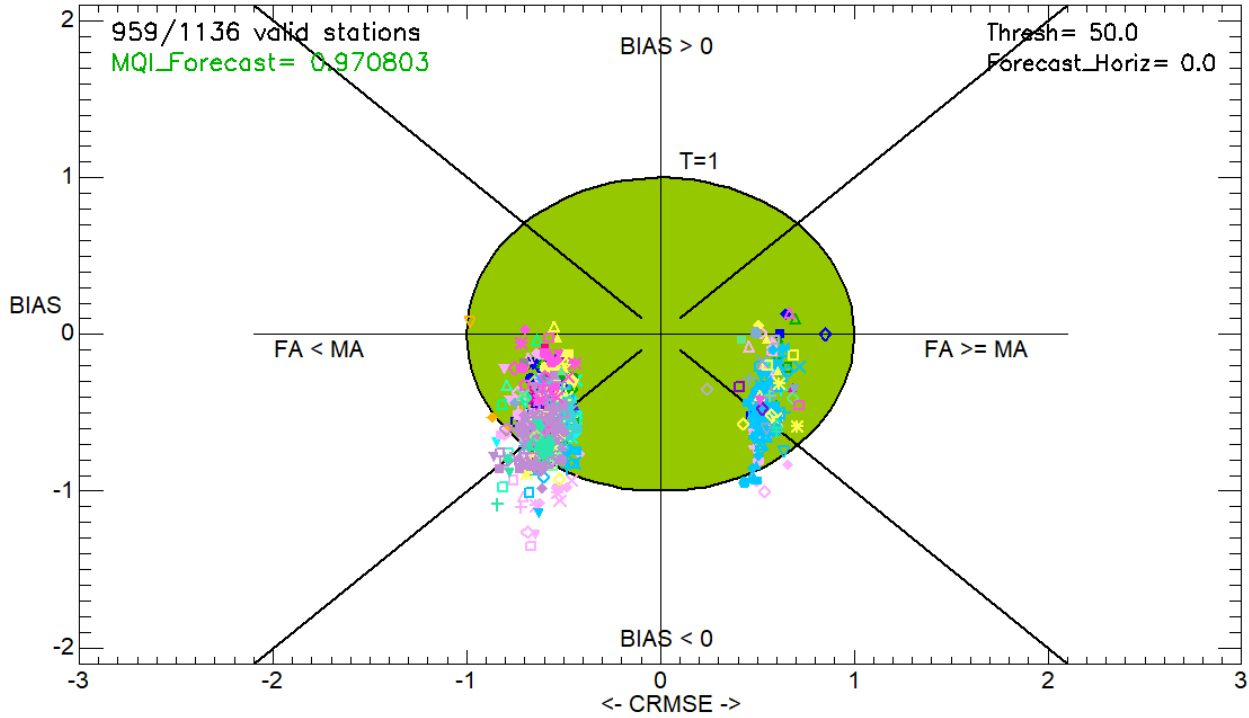
Standard assessment MOQ outcomes (i.e., MQI value associated to the 90th percentile worst station) for all available pollutants are summarized in Table C1 for all forecast applications.

Table C1 Standard assessment MQI values (associated to the 90th percentile worst station) for all forecast applications.

	<u>NO₂</u>	<u>O₃</u>	<u>PM10</u>	<u>PM2.5</u>
<u>FA1</u>	<u>0.865</u>	<u>0.619</u>	<u>1.267</u>	<u>0.776</u>
<u>FA2</u>	<u>0.831</u>	<u>0.698</u>	<u>0.941</u>	<u>0.700</u>
<u>FA3</u>			<u>0.479</u>	
<u>FA4</u>	<u>1.009</u>	<u>0.696</u>	<u>0.943</u>	<u>1.009</u>
<u>FA5</u>	<u>0.685</u>	<u>0.570</u>		<u>0.528</u>

570

FORECAST TARGET PLOT (OU) PM10



○ AL0204A	+ AT10003	■ AT30101	▲ AT31401	× AT32701	Strt/end Ind: 1-8760 Model (s): v51 Parameter: PM10 Scen: 2018 Extra Values: 50/0 Season: Year Day hours: All 24h Time Average: Preserve Daily stats: Mean
○ AT0ENK1	● AT10035	▼ AT30201	▼ AT31413	○ AT4S108	
○ AT0ILL1	■ AT11007	▼ AT30301	● AT31703	○ AT4S125	
▲ AT0KLLH1	× AT2KA71	+ AT30302	■ AT31901	▲ AT4S156	
▼ AT0PIL1	○ AT2M226	○ AT30401	■ AT31904	▼ AT4S184	
● AT0VOR1	□ AT2SP10	■ AT30502	▼ AT32301	● AT4S404	
■ AT0ZOE2	▲ AT2SP18	■ AT30603	+ AT32401	■ AT4S406	
▲ AT10001	▼ AT2WO25	○ AT30701	● AT32501	▲ AT4S407	
▼ AT10002	● AT2WO35	■ AT31301	■ AT32604	▼ AT4S409	

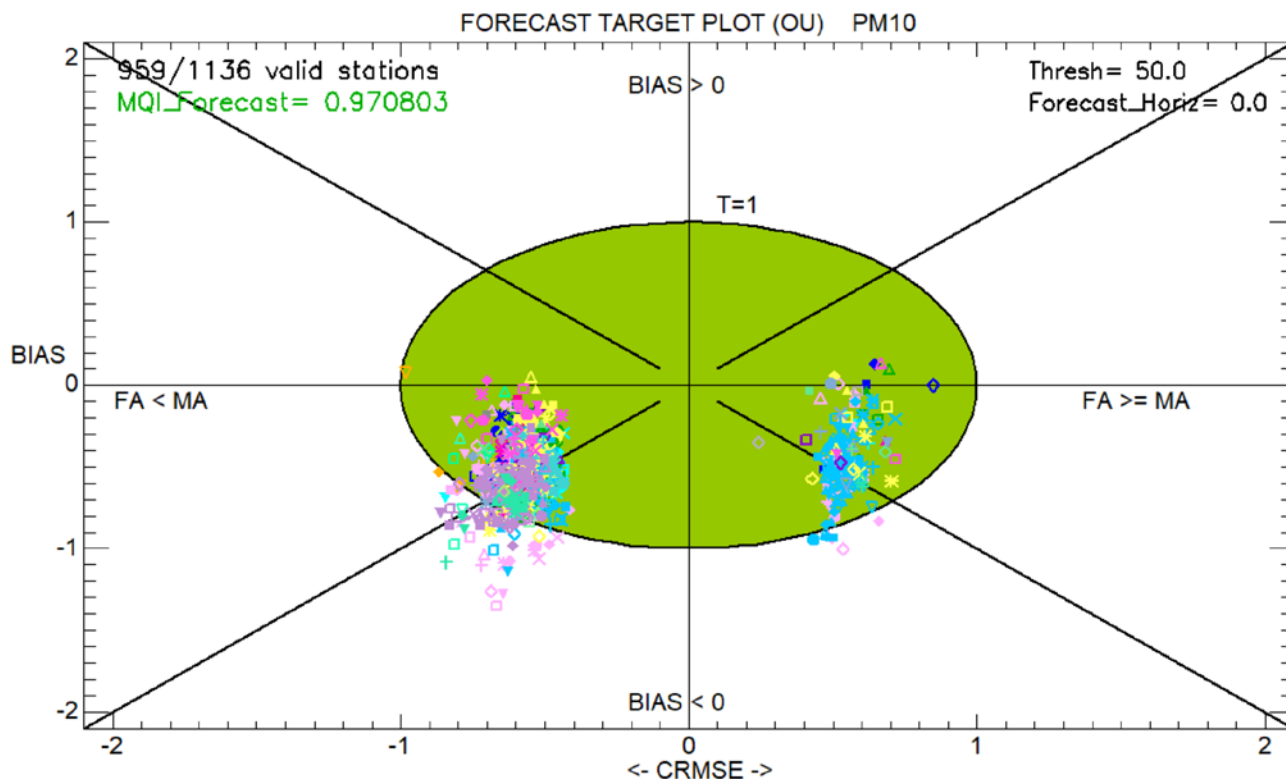


Figure C4D1: FA1 Forecast Target Plot for PM10, removing Turkish monitoring stations from the validation data set.

575 **Appendix E**

580 The effect on MQI_f outcomes of lowering measurement uncertainty estimates is investigated here. More in detail, the values of $U_f(RV)$ parameters in Table A1 (i.e. the estimates of the relative uncertainty around the reference value, defining the asymptotic behaviour of the functions of Figure A1) were reduced by 25% and 50% for all the pollutants and the MQI_f were recalculated for the different forecast applications. Figure E1 shows the results for all available data: FA1, FA2, FA4 outcomes for the current forecast day (all pollutants available), and FA3 outcomes along a three-days forecast horizon (only PM10 available). Different colours refer to results based on different $U_f(RV)$ values: 1 $U_f(RV)$ indicates the original values in Table A1, 0.75 $U_f(RV)$ and 0.50 $U_f(RV)$ refer to 25% and 50% reductions, respectively. Indeed, the 50% reduction decreases the $U_f(RV)$ values to 0.12 (NO_2), 0.09 (O_3), 0.14 (PM10), 0.18 (PM2.5), i.e. well below the data quality objective values set by the current European legislation (European Union, 2008), namely 15%, 15% and 25% for NO_2 , O_3 , and particulate matter.

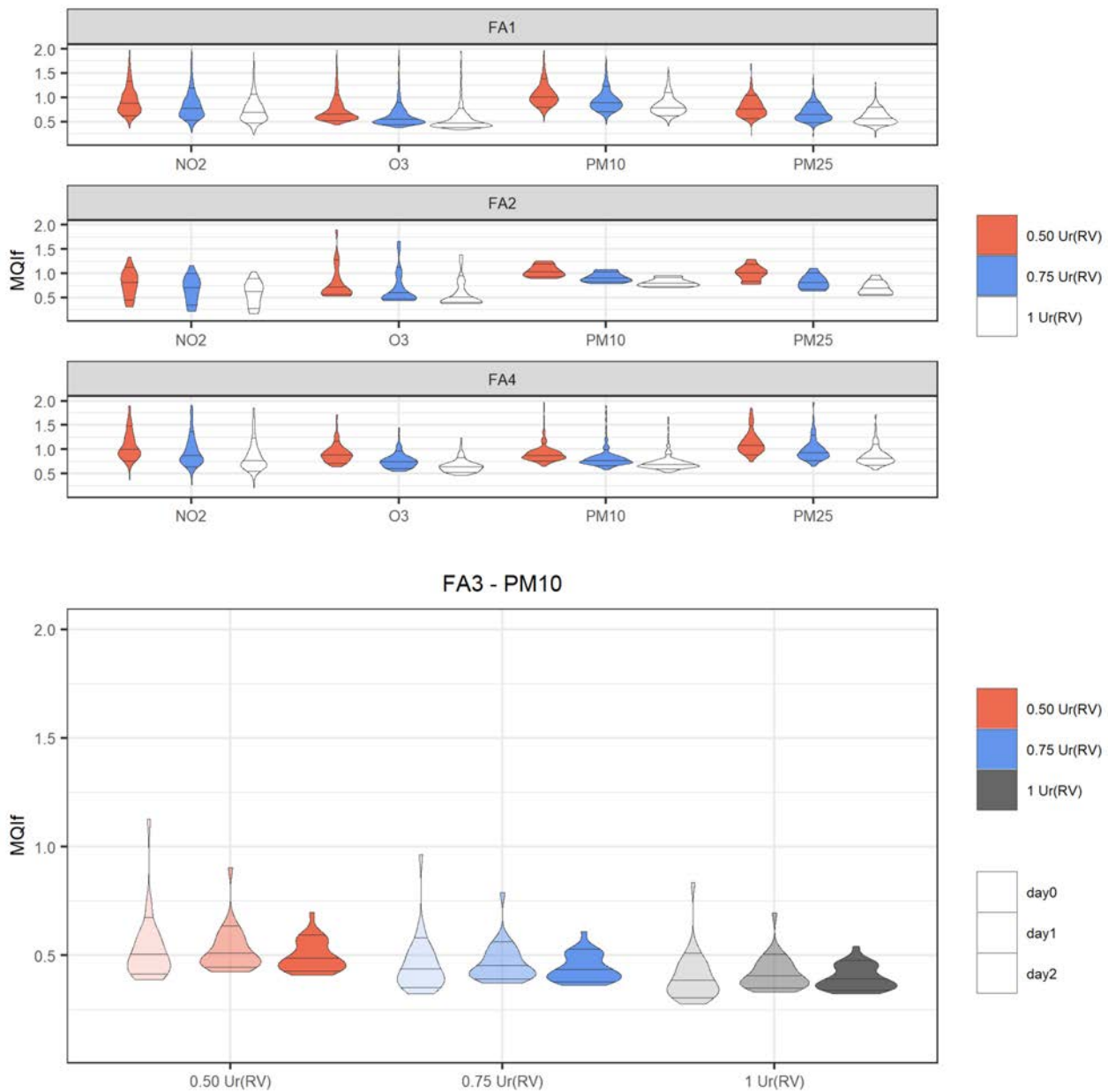
585 The results of the sensitivity analysis are provided by means of the violin plots (Hintze and Nelson, 1998), showing the distributions of the MQI_f values computed for each monitoring station. In other words, each violin refers to all the data

provided within the corresponding Forecast Target Plot, giving in a single plot an overall view of all the outcomes available. Three lines were added to the display of each violin, indicating the 10th, the 50th and 90th percentiles of the distributions.

590 Results show that both MOI_f values and the shape of their distribution depend on both the forecast application and the pollutant. Within this context, changing $U_r(RV)$ values induces a very slight effect on the shape of the MOI_f values distribution, apart from the case of PM2.5 for FA2, where a small amount of data is available (11 monitoring stations). On the contrary, as expected, changing $U_r(RV)$ values turns out in variations of MOI_f values, which get higher as $U_r(RV)$ gets lower, to a different extent depending on the forecast application and the pollutant. Generally, variations tend to be lower if

595 data availability is higher. Concerning the main MQO_f criterion fulfilment (i.e. the 90th percentile of the MOI_f values is lower than 1), being based on a categorical answer (yes/no), it changes or not mainly depending on the performances of the reference analysis (1 $U_r(RV)$). The same answer is maintained both in case of very good performances (MOI_f 90th percentile value largely lower than 1) and in case the criterion is not fulfilled even in the reference analysis (MOI_f 90th percentile value already higher than 1). When MOI_f 90th percentile value is lower but quite close to 1, MQO_f criterion fulfilment is of course

600 more sensitive to measurement uncertainty estimate. Indeed, this is expected and it is a typical shortcoming of the usage of criteria based on categorical answers.



605

Figure E1: Effect on the distribution of MOI_f values of lowering $U_f(RV)$. Top: FA1, FA2, FA4 outcomes for the current forecast day (all the pollutants available). Bottom: FA3 outcomes along a three-days forecast horizon (only PM10 available).

Code and data availability. The DELTA Tool software and all datasets generated and analysed during the current study are available on Zenodo at <https://doi.org/10.5281/zenodo.7949868>. DELTA Tool software is available at

610 ~~<https://aqm.jrc.ec.europa.eu/Section/Assessment/Download> and can be downloaded after free registration. The datasets generated and analysed during the current study are available from the corresponding author on request.~~

Author contributions. KC, PD, SJ, AM, AP, PT and LV contributed to the study conception and design. PT conceived part of the methodology and KC developed the software. Material preparation, data collection and analysis were performed by
615 MA, RA, AB, AD, CG, GG, AM, TP, MS, LV and SV. The first draft of the manuscript was written by AM, AP and LV and all authors commented on previous versions of the manuscript. All authors read and approved the final manuscript.

Competing interests. The authors declare that they have no conflict of interest.

620 **Acknowledgements.** Thanks are due to FCT/MCTES for the contract granted to Carla Gama (2021.00732.CEECIND) and for the financial support to the CESAM Associated Laboratory (UIDB/50017/2020 + UIDP/50017/2020).

The analysis for the Irish air quality forecasts has been supported by the LIFE Emerald project “Emissions Modelling and Forecasting of Air in Ireland” which is co-funded by the European Union, under grant agreement No. LIFE19
625 GIE/IE/001101. VITO thanks the Irish EPA for sharing all observation data and continuous valuable feedback on the model developments.

NINFA simulation over Po Valley and Slovenia was developed under the project LIFE-IP PREPAIR (Po Regions Engaged to Policies of AIR), which was co-funded by the European Union LIFE Program, in 2016, Grant Number LIFE15
630 IPE/IT/000013. Acknowledgments are given to all the Beneficiaries of the project LIFE-IP PREPAIR: Emilia-Romagna Region (Project Coordinator), Veneto Region, Lombardy Region, Piedmont Region, Friuli Venezia Giulia Region, Autonomous Province of Trento, Regional Agency for Environment of Emilia-Romagna, Regional Agency for Environment of Veneto, Regional Agency for Environment of Piedmont, Regional Agency for Environmental Protection of Lombardy, Environmental Protection Agency of Valle d’Aosta, Environmental Protection Agency of Friuli Venezia Giulia, Slovenian Environment Agency, Municipality of Bologna, Municipality of Milan, City of Turin, ART-ER, Lombardy Foundation for Environment

635 The forecasts for Kosovo were developed under the project: “Supply of project management, air quality information management, behaviour change and communication services” managed by the Millennium Foundation Kosovo (MFK), funded by the Millennium Challenge Corporation (MCC). Acknowledgments are given to all the Beneficiaries and Participants of the project for Kosovo, including Millennium Foundation Kosovo; the Hydrometeorological Institute of Kosovo and the National Institute of Public Health of Kosovo.

640 References

- Adani, M., Piersanti, A., Ciancarella, L., D'Isidoro, M., Villani, M. G., and Vitali, L.: Preliminary Tests on the Sensitivity of the FORAIR_IT Air Quality Forecasting System to Different Meteorological Drivers, *Atmosphere*, 11, 574, <https://doi.org/10.3390/atmos11060574>, 2020.
- 645 Adani, M., D'Isidoro, M., Mircea, M., Guarnieri, G., Vitali, L., D'Elia, I., Ciancarella, L., Gualtieri, M., Briganti, G., Cappelletti, A., Piersanti, A., Stracquadanio, M., Righini, G., Russo, F., Cremona, G., Villani, M. G., and Zanini, G.: Evaluation of air quality forecasting system FORAIR-IT over Europe and Italy at high resolution for year 2017, *Atmospheric Pollut. Res.*, 13, 101456, <https://doi.org/10.1016/j.apr.2022.101456>, 2022.
- 650 Agarwal, S., Sharma, S., R., S., Rahman, M. H., Vranckx, S., Maiheu, B., Blyth, L., Janssen, S., Gargava, P., Shukla, V. K., and Batra, S.: Air quality forecasting using artificial neural networks with real time dynamic error correction in highly polluted regions, *Sci. Total Environ.*, 735, 139454, <https://doi.org/10.1016/j.scitotenv.2020.139454>, 2020.
- Alfaro, S. C. and Gomes, L.: Modeling mineral aerosol production by wind erosion: Emission intensities and aerosol size distributions in source areas, *J. Geophys. Res. Atmospheres*, 106, 18075–18084, <https://doi.org/10.1029/2000JD900339>, 2001.
- 655 Bai, L., Wang, J., Ma, X., and Lu, H.: Air Pollution Forecasts: An Overview, *Int. J. Environ. Res. Public Health*, 15, 780, <https://doi.org/10.3390/ijerph15040780>, 2018.
- Baklanov, A. and Zhang, Y.: Advances in air quality modeling and forecasting, *Glob. Transit.*, 2, 261–270, <https://doi.org/10.1016/j.glt.2020.11.001>, 2020.
- 660 Baklanov, A., Schlünzen, K., Suppan, P., Baldasano, J., Brunner, D., Aksoyoglu, S., Carmichael, G., Douros, J., Flemming, J., Forkel, R., Galmarini, S., Gauss, M., Grell, G., Hirtl, M., Joffre, S., Jorba, O., Kaas, E., Kaasik, M., Kallos, G., Kong, X., Korsholm, U., Kurganskiy, A., Kushta, J., Lohmann, U., Mahura, A., Manders-Groot, A., Maurizi, A., Moussiopoulos, N., Rao, S. T., Savage, N., Seigneur, C., Sokhi, R. S., Solazzo, E., Solomos, S., Sørensen, B., Tsegas, G., Vignati, E., Vogel, B., and Zhang, Y.: Online coupled regional meteorology chemistry models in Europe: current status and prospects, *Atmospheric Chem. Phys.*, 14, 317–398, <https://doi.org/10.5194/acp-14-317-2014>, 2014.
- 665 Baldauf, M., Seifert, A., Förstner, J., Majewski, D., Raschendorfer, M., and Reinhardt, T.: Operational Convective-Scale Numerical Weather Prediction with the COSMO Model: Description and Sensitivities, *Mon. Weather Rev.*, 139, 3887–3905, <https://doi.org/10.1175/MWR-D-10-05013.1>, 2011.
- Borrego, C., Monteiro, A., Ferreira, J., Miranda, A. I., Costa, A. M., Carvalho, A. C., and Lopes, M.: Procedures for estimation of modelling uncertainty in air quality assessment, *Environ. Int.*, 34, 613–620, <https://doi.org/10.1016/j.envint.2007.12.005>, 2008.
- 670 Boylan, J. W. and Russell, A. G.: PM and light extinction model performance metrics, goals, and criteria for three-dimensional air quality models, *Atmos. Environ.*, 40, 4946–4959, <https://doi.org/10.1016/j.atmosenv.2005.09.087>, 2006.
- Cabaneros, S. M., Calautit, J. K., and Hughes, B. R.: A review of artificial neural network models for ambient air pollution prediction, *Environ. Model. Softw.*, 119, 285–304, <https://doi.org/10.1016/j.envsoft.2019.06.014>, 2019.
- 675 Carnevale, C., Finzi, G., Pederzoli, A., Pisoni, E., Thunis, P., Turrini, E., and Volta, M.: A methodology for the evaluation of re-analyzed PM10 concentration fields: a case study over the PO Valley, *Air Qual. Atmosphere Health*, 8, 533–544, <https://doi.org/10.1007/s11869-014-0307-2>, 2015.

- Chang, J. C. and Hanna, S. R.: Air quality model performance evaluation, *Meteorol. Atmospheric Phys.*, 87, 167–196, <https://doi.org/10.1007/s00703-003-0070-7>, 2004.
- 680 Chemel, C., Sokhi, R. S., Yu, Y., Hayman, G. D., Vincent, K. J., Dore, A. J., Tang, Y. S., Prain, H. D., and Fisher, B. E. A.: Evaluation of a CMAQ simulation at high resolution over the UK for the calendar year 2003, *Atmos. Environ.*, 44, 2927–2939, <https://doi.org/10.1016/j.atmosenv.2010.03.029>, 2010.
- D’Elia, I., Briganti, G., Vitali, L., Piersanti, A., Righini, G., D’Isidoro, M., Cappelletti, A., Mircea, M., Adani, M., Zanini, G., and Ciancarella, L.: Measured and modelled air quality trends in Italy over the period 2003–2010, *Atmospheric Chem. Phys.*, 21, 10825–10849, <https://doi.org/10.5194/acp-21-10825-2021>, 2021.
- 685 Dennis, R., Fox, T., Fuentes, M., Gilliland, A., Hanna, S., Hogrefe, C., Irwin, J., Rao, S. T., Scheffe, R., Schere, K., Steyn, D., and Venkatram, A.: A framework for evaluating regional-scale numerical photochemical modeling systems, *Environ. Fluid Mech.*, 10, 471–489, <https://doi.org/10.1007/s10652-009-9163-2>, 2010.
- Doms, G. and Baldauf, M.: A Description of the Non Hydrostatic Regional COSMO-Model. Part I: Dynamics and Numeric. User Guide Documentation. Available online: www.cosmo-model.org, 2018.
- 690 Eder, B., Kang, D., Rao, S. T., Mathur, R., Yu, S., Otte, T., Schere, K., Wayland, R., Jackson, S., Davidson, P., McQueen, J., and Bridgers, G.: Using National Air Quality Forecast Guidance to Develop Local Air Quality Index Forecasts, *Bull. Am. Meteorol. Soc.*, 91, 313–326, <https://doi.org/10.1175/2009BAMS2734.1>, 2010.
- Emery, C., Liu, Z., Russell, A. G., Odman, M. T., Yarwood, G., and Kumar, N.: Recommendations on statistics and benchmarks to assess photochemical model performance, *J. Air Waste Manag. Assoc.*, 67, 582–598, <https://doi.org/10.1080/10962247.2016.1265027>, 2017.
- 695 EPA: Guidelines for Developing an Air Quality (Ozone and PM_{2.5}) Forecasting Program, EPA-456/R-03-002 June 2003, 2003.
- European Union: Directive 2004/107/EC of the European Parliament and of the Council of 15 December 2004 relating to arsenic, cadmium, mercury, nickel and polycyclic aromatic hydrocarbons in ambient air, 2004.
- 700 European Union: Directive 2008/50/EC of the European Parliament and of the Council of 21 May 2008 on ambient air quality and cleaner air for Europe, OJ L, 152, 2008.
- European Union: Proposal for a DIRECTIVE OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL on ambient air quality and cleaner air for Europe (recast), 2022.
- 705 Georgieva, E., Syrakov, D., Prodanova, M., Etopolska, I., and Slavov, K.: Evaluating the performance of WRF-CMAQ air quality modelling system in Bulgaria by means of the DELTA tool, *Int. J. Environ. Pollut.*, 57, 272–284, <https://doi.org/10.1504/IJEP.2015.074512>, 2015.
- Ginoux, P., Chin, M., Tegen, I., Prospero, J. M., Holben, B., Dubovik, O., and Lin, S.-J.: Sources and distributions of dust aerosols simulated with the GOCART model, *J. Geophys. Res. Atmospheres*, 106, 20255–20273, <https://doi.org/10.1029/2000JD000053>, 2001.
- 710 Guenther, A., Karl, T., Harley, P., Wiedinmyer, C., Palmer, P. I., and Geron, C.: Estimates of global terrestrial isoprene emissions using MEGAN (Model of Emissions of Gases and Aerosols from Nature), *Atmospheric Chem. Phys.*, 6, 3181–3210, <https://doi.org/10.5194/acp-6-3181-2006>, 2006.

- Hanna, S. R. and Chang, J.: Setting Acceptance Criteria for Air Quality Models, in: Air Pollution Modeling and its Application XXI, Dordrecht, 479–484, https://doi.org/10.1007/978-94-007-1359-8_80, 2012.
- 715 Hintze, J. L. and Nelson, R. D.: Violin Plots: A Box Plot-Density Trace Synergism, *Am. Stat.*, *52*, 181–184, <https://doi.org/10.1080/00031305.1998.10480559>, 1998.
- Hooyberghs, J., Mensink, C., Dumont, G., Fierens, F., and Brasseur, O.: A neural network forecast for daily average PM10 concentrations in Belgium, *Atmos. Environ.*, *39*, 3279–3289, <https://doi.org/10.1016/j.atmosenv.2005.01.050>, 2005.
- 720 Janssen, S. and Thunis, P.: FAIRMODE Guidance Document on Modelling Quality Objectives and Benchmarking (version 3.3), EUR 31068 EN, Publications Office of the European Union, Luxembourg, 2022, ISBN 978-92-76-52425-0, <https://doi.org/10.2760/41988>, JRC129254, 2022.
- Janssen, S., Dumont, G., Fierens, F., and Mensink, C.: Spatial interpolation of air pollution measurements using CORINE land cover data, *Atmos. Environ.*, *42*, 4884–4903, <https://doi.org/10.1016/j.atmosenv.2008.02.043>, 2008.
- 725 Kang, D., Eder, B. K., Stein, A. F., Grell, G. A., Peckham, S. E., and McHenry, J.: The New England Air Quality Forecasting Pilot Program: Development of an Evaluation Protocol and Performance Benchmark, *J. Air Waste Manag. Assoc.*, *55*, 1782–1796, <https://doi.org/10.1080/10473289.2005.10464775>, 2005.
- Kingma, D. P. and Ba, J.: Adam: A Method for Stochastic Optimization. rXiv preprint arXiv:1412.6980. <https://doi.org/10.48550/arXiv.1412.6980>, 2014.
- 730 Knaff, J. A. and Landsea, C. W.: An El Niño–Southern Oscillation Climatology and Persistence (CLIPER) Forecasting Scheme, *Weather Forecast.*, *12*, 633–652, [https://doi.org/10.1175/1520-0434\(1997\)012<0633:AENOSO>2.0.CO;2](https://doi.org/10.1175/1520-0434(1997)012<0633:AENOSO>2.0.CO;2), 1997.
- Kukkonen, J., Olsson, T., Schultz, D. M., Baklanov, A., Klein, T., Miranda, A. I., Monteiro, A., Hirtl, M., Tarvainen, V., Boy, M., Peuch, V.-H., Poupkou, A., Kioutsioukis, I., Finardi, S., Sofiev, M., Sokhi, R., Lehtinen, K. E. J., Karatzas, K., San José, R., Astitha, M., Kallos, G., Schaap, M., Reimer, E., Jakobs, H., and Eben, K.: A review of operational, regional-scale, chemical weather forecasting models in Europe, *Atmospheric Chem. Phys.*, *12*, 1–87, <https://doi.org/10.5194/acp-12-1-2012>, 2012.
- 735 Kushta, J., Georgiou, G. K., Proestos, Y., Christoudias, T., Thunis, P., Savvides, C., Papadopoulos, C., and Lelieveld, J.: Evaluation of EU air quality standards through modeling and the FAIRMODE benchmarking methodology, *Air Qual. Atmosphere Health*, *12*, 73–86, <https://doi.org/10.1007/s11869-018-0631-z>, 2019.
- 740 Mailler, S., Menut, L., Khvorostyanov, D., Valari, M., Couvidat, F., Siour, G., Turquety, S., Briant, R., Tuccella, P., Bessagnet, B., Colette, A., Létinois, L., Markakis, K., and Meleux, F.: CHIMERE-2017: from urban to hemispheric chemistry-transport modeling, *Geosci. Model Dev.*, *10*, 2397–2423, <https://doi.org/10.5194/gmd-10-2397-2017>, 2017.
- 745 Marécal, V., Peuch, V.-H., Andersson, C., Andersson, S., Arteta, J., Beekmann, M., Benedictow, A., Bergström, R., Bessagnet, B., Cansado, A., Chérour, F., Colette, A., Coman, A., Curier, R. L., Denier van der Gon, H. a. C., Drouin, A., Elbern, H., Emili, E., Engelen, R. J., Eskes, H. J., Foret, G., Friese, E., Gauss, M., Giannaros, C., Guth, J., Joly, M., Jaumouillé, E., Josse, B., Kadyrov, N., Kaiser, J. W., Krajsek, K., Kuenen, J., Kumar, U., Liora, N., Lopez, E., Malherbe, L., Martinez, I., Melas, D., Meleux, F., Menut, L., Moinat, P., Morales, T., Parmentier, J., Piacentini, A., Plu, M., Poupkou, A., Queguiner, S., Robertson, L., Rouil, L., Schaap, M., Segers, A., Sofiev, M., Tarasson, L., Thomas, M., Timmermans, R., Valdebenito, Á., van Velthoven, P., van Versendaal, R., Vira, J., and Ung, A.: A regional air quality forecasting system over Europe: the MACC-II daily ensemble production, *Geosci. Model Dev.*, *8*, 2777–2813, <https://doi.org/10.5194/gmd-8-2777-2015>, 2015.
- 750

- Marongiu, A., Angelino, E., Moretti, M., Malvestiti, G., and Fossati, G.: Atmospheric Emission Sources in the Po-Basin from the LIFE-IP PREPAIR Project, *Open J. Air Pollut.*, 11, 70–83, <https://doi.org/10.4236/ojap.2022.113006>, 2022.
- 755 Masood, A. and Ahmad, K.: A review on emerging artificial intelligence (AI) techniques for air pollution forecasting: Fundamentals, application and performance, *J. Clean. Prod.*, 322, 129072, <https://doi.org/10.1016/j.jclepro.2021.129072>, 2021.
- Menut, L., Bessagnet, B., Khvorostyanov, D., Beekmann, M., Blond, N., Colette, A., Coll, I., Curci, G., Foret, G., Hodzic, A., Mailler, S., Meleux, F., Monge, J.-L., Pison, I., Siour, G., Turquety, S., Valari, M., Vautard, R., and Vivanco, M. G.: CHIMERE 2013: a model for regional atmospheric composition modelling, *Geosci. Model Dev.*, 6, 981–1028, <https://doi.org/10.5194/gmd-6-981-2013>, 2013.
- 760 Mircea, M., Ciancarella, L., Briganti, G., Calori, G., Cappelletti, A., Cionni, I., Costa, M., Cremona, G., D’Isidoro, M., Finardi, S., Pace, G., Piersanti, A., Righini, G., Silibello, C., Vitali, L., and Zanini, G.: Assessment of the AMS-MINNI system capabilities to simulate air quality over Italy for the calendar year 2005, *Atmos. Environ.*, 84, 178–188, <https://doi.org/10.1016/j.atmosenv.2013.11.006>, 2014.
- 765 Mittermaier, M. P.: The Potential Impact of Using Persistence as a Reference Forecast on Perceived Forecast Skill, *Weather Forecast.*, 23, 1022–1031, <https://doi.org/10.1175/2008WAF2007037.1>, 2008.
- Monteiro, A., Lopes, M., Miranda, A. I., Borrego, C., and Robert Vautard: Air pollution forecast in Portugal: a demand from the new air quality framework directive, *Int. J. Environ. Pollut.*, 25, 4–15, <https://doi.org/10.1504/IJEP.2005.007650>, 2005.
- Monteiro, A., Miranda, A. I., Borrego, C., and Vautard, R.: Air quality assessment for Portugal, *Sci. Total Environ.*, 373, 22–31, <https://doi.org/10.1016/j.scitotenv.2006.10.014>, 2007a.
- 770 Monteiro, A., Miranda, A. I., Borrego, C., Vautard, R., Ferreira, J., and Perez, A. T.: Long-term assessment of particulate matter using CHIMERE model, *Atmos. Environ.*, 41, 7726–7738, <https://doi.org/10.1016/j.atmosenv.2007.06.008>, 2007b.
- 775 Monteiro, A., Durka, P., Flandorfer, C., Georgieva, E., Guerreiro, C., Kushta, J., Malherbe, L., Maiheu, B., Miranda, A. I., Santos, G., Stocker, J., Trimpeeneers, E., Tognet, F., Stortini, M., Wesseling, J., Janssen, S., and Thunis, P.: Strengths and weaknesses of the FAIRMODE benchmarking methodology for the evaluation of air quality models, *Air Qual. Atmosphere Health*, 11, 373–383, <https://doi.org/10.1007/s11869-018-0554-8>, 2018.
- [Olesen, H. R.: Toward the Establishment of a Common Framework for Model Evaluation, in: Air Pollution Modeling and Its Application XI, edited by: Gryning, S.-E. and Schiermeier, F. A., Springer US, Boston, MA, 519–528, https://doi.org/10.1007/978-1-4615-5841-5_54, 1996.](https://doi.org/10.1007/978-1-4615-5841-5_54)
- 780 Pernigotti, D., Gerboles, M., Belis, C. A., and Thunis, P.: Model quality objectives based on measurement uncertainty. Part II: NO₂ and PM₁₀, *Atmos. Environ.*, 79, 869–878, <https://doi.org/10.1016/j.atmosenv.2013.07.045>, 2013.
- Raffaelli, K., Deserti, M., Stortini, M., Amorati, R., Vasconi, M., and Giovannini, G.: Improving Air Quality in the Po Valley, Italy: Some Results by the LIFE-IP-PREPAIR Project, *Atmosphere*, 11, 429, <https://doi.org/10.3390/atmos11040429>, 2020.
- 785 Rahman, M. H., Agarwal, S., Sharma, S., Suresh, R., Kundu, S., Vranckx, S., Maiheu, B., Blyth, L., Janssen, S., Jorge, S., Gargava, P., Shukla, V. K., and Batra, S.: High-Resolution Mapping of Air Pollution in Delhi Using Detrended Kriging Model, *Environ. Model. Assess.*, 28, 39–54, <https://doi.org/10.1007/s10666-022-09842-5>, 2023.

- Russell, A. and Dennis, R.: NARSTO critical review of photochemical models and modeling, *Atmos. Environ.*, 34, 2283–2324, [https://doi.org/10.1016/S1352-2310\(99\)00468-9](https://doi.org/10.1016/S1352-2310(99)00468-9), 2000.
- Ryan, W. F.: The air quality forecast rote: Recent changes and future challenges, *J. Air Waste Manag. Assoc.*, 66, 576–596, <https://doi.org/10.1080/10962247.2016.1151469>, 2016.
- Seigneur, C., Pun, B., Pai, P., Louis, J.-F., Solomon, P., Emery, C., Morris, R., Zahniser, M., Worsnop, D., Koutrakis, P., White, W., and Tombach, I.: Guidance for the Performance Evaluation of Three-Dimensional Air Quality Modeling Systems for Particulate Matter and Visibility, *J. Air Waste Manag. Assoc.*, 50, 588–599, <https://doi.org/10.1080/10473289.2000.10464036>, 2000.
- 795 Skamarock, C., Klemp, B., Dudhia, J., Gill, O., Barker, D., Duda, G., Huang, X., Wang, W., and Powers, G.: A Description of the Advanced Research WRF Version 3. NCAR Tech. Note NCAR/TN-475+STR, <https://doi.org/10.5065/D68S4MVH>, 2008.
- 800 Sokhi, R. S., Moussiopoulos, N., Baklanov, A., Bartzis, J., Coll, I., Finardi, S., Friedrich, R., Geels, C., Grönholm, T., Halenka, T., Ketzler, M., Maragkidou, A., Matthias, V., Moldanova, J., Ntziachristos, L., Schäfer, K., Suppan, P., Tsegas, G., Carmichael, G., Franco, V., Hanna, S., Jalkanen, J.-P., Velders, G. J. M., and Kukkonen, J.: Advances in air quality research – current and emerging challenges, *Atmospheric Chem. Phys.*, 22, 4615–4703, <https://doi.org/10.5194/acp-22-4615-2022>, 2022.
- Stortini, M., Arvani, B., and Deserti, M.: Operational Forecast and Daily Assessment of the Air Quality in Italy: A Copernicus-CAMS Downstream Service, *Atmosphere*, 11, 447, <https://doi.org/10.3390/atmos11050447>, 2020.
- 805 Szopa, S., Foret, G., Menut, L., and Cozic, A.: Impact of large scale circulation on European summer surface ozone and consequences for modelling forecast, *Atmos. Environ.*, 43, 1189–1195, <https://doi.org/10.1016/j.atmosenv.2008.10.039>, 2009.
- Tesche, T. W., Lurmann, F. R., Roth, P. M., Georgopoulos, P., and Seinfeld, J. H.: Improvement of procedures for evaluating photochemical models. Final report, Radian Corp., Sacramento, CA (USA), 1990.
- 810 Thunis, P., Georgieva, E., and Pederzoli, A.: A tool to evaluate air quality model performances in regulatory applications, *Environ. Model. Softw.*, 38, 220–230, <https://doi.org/10.1016/j.envsoft.2012.06.005>, 2012a.
- Thunis, P., Pederzoli, A., and Pernigotti, D.: Performance criteria to evaluate air quality modeling applications, *Atmos. Environ.*, 59, 476–482, <https://doi.org/10.1016/j.atmosenv.2012.05.043>, 2012b.
- 815 Thunis, P., Pernigotti, D., and Gerboles, M.: Model quality objectives based on measurement uncertainty. Part I: Ozone, *Atmos. Environ.*, 79, 861–868, <https://doi.org/10.1016/j.atmosenv.2013.05.018>, 2013.
- Zhang, B., Rong, Y., Yong, R., Qin, D., Li, M., Zou, G., and Pan, J.: Deep learning for air pollutant concentration prediction: A review, *Atmos. Environ.*, 290, 119347, <https://doi.org/10.1016/j.atmosenv.2022.119347>, 2022.
- Zhang, Y., Bocquet, M., Mallet, V., Seigneur, C., and Baklanov, A.: Real-time air quality forecasting, part I: History, techniques, and current status, *Atmos. Environ.*, 60, 632–655, <https://doi.org/10.1016/j.atmosenv.2012.06.031>, 2012.