



How accurate are operational dust models in predicting Particulate Matter (PM) levels in the Eastern Mediterranean Region? Insights from PM Surface Concentrations

Andreas Eleftheriou¹, Petros Mouzourides¹, Panayiotis Kouis², Nikos Kalivitis³, Itzhak Katra⁴, Emily Vasiliadou⁵, Chrysanthos Savvides⁵, Panayiotis Yiallourous², Marina K.-A. Neophytou¹

¹Environmental Fluid Mechanics Laboratory, Department of Civil and Environmental Engineering, University of Cyprus, Nicosia, Cyprus

²Medical School, Cyprus University, Cyprus

³Department of Chemistry, University of Crete, Greece

⁴Department of Environmental, Geoinformatics, and Urban Planning Sciences, Ben-Gurion University of the Negev, Israel

⁵Department of Labour Inspection, Ministry of Labour and Social Insurance, Cyprus

Correspondence to: Petros Mouzourides (pmouzou@ucy.ac.cy)

Abstract. This study provides the first comprehensive assessment of eleven operational dust forecast models and a multi-model ensemble in predicting ground-level Particulate Matter (PM) concentrations in the Eastern Mediterranean Region (EMR), with a focus on Cyprus, Greece, and Israel. Ground-based observations from regional background stations support model performance assessment across different PM fractions (PM₁₀, PM_{2.5}, and coarse particles), using established statistical metrics (correlation coefficient, R, Mean Bias, MB, and Root Mean Square Error, RMSE). The results reveal substantial variability in accuracy, with R values ranging from −0.24 to 0.91 depending on site and event subset. NASA-GEOS consistently achieves the highest correlation (R = 0.71 at Cyprus), indicating accurate representation of dust transport. In contrast, SILAM and EMA-REG4 perform poorly, with low correlations (R = 0.10 and −0.24, respectively) and significant estimation errors (MB = −90.34 µg/m³ for EMA-REG4). The NOAA-WRF model effectively captures extreme dust events, with R = 0.91 during the 95th percentile of PM concentrations in Greece. Most models perform better for coarse PM, with the BOOT methodology indicating reduced scatter and bias during dust storm days. However, no model performs optimally across all sites and conditions, highlighting the need for location-specific tuning and evaluation. The study underscores the importance of refining model configurations and improving parameterizations to enhance forecast accuracy. Future efforts should incorporate localized data and further develop region-specific models to improve the operational use of these systems in early warning protocols for mitigating public health impacts.

1. Introduction

Desert Dust Storm (DDS) events have been documented since ancient times, with early references found in texts such as “Histories” by Herodotus (Herodotus III: 86-88). These events impact human health and infrastructure, yet they remain understudied until recent decades.

Over the past few years, research into DDSs has expanded rapidly. Recognising their global impacts, the United Nations General Assembly designates 2025-2034 as the Decade for combatting DDSs (United Nations, 2024). This growing attention reflects the mounting evidence linking DDSs to respiratory and cardiovascular illnesses (Lorentzou et al., 2019) and mental health effects (Jones, 2023), regardless of the source region (Lwin et al., 2023). In parallel, climate change, population growth, and geopolitical instability increasingly exacerbate desertification, driving more frequent and intense DDSs and reinforcing transboundary dust



transport patterns (Eleftheriou et al., 2023). Therefore, monitoring both the source and downwind impacts of DDSs remains critical for assessing and mitigating their effects.

To address air quality concerns, European directives (e.g. 2004/107/EC, 2008/50/EC, and 2015/1480/EU) and the World Health Organization (WHO) air quality guidelines (2021), establish threshold values to regulate pollution from both anthropogenic and natural sources. These regulatory frameworks enable authorities to track elevated PM concentrations from human activities, while also capturing high PM levels during natural events such as DDSs. During such episodes, PM concentrations often exceed 150–200 $\mu\text{g}/\text{m}^3$, well above regulatory limits. This capability supports timely public warnings and facilitates health-protective measures when air quality deteriorates due to dust intrusions. The integration of monitoring networks with operational forecast models, such as METAL-WRF (Solomos et al., 2023), enhances our understanding of DDS dynamics and supports compliance with air quality regulations. These systems also strengthen community resilience, providing predictive capacity to mitigate both anthropogenic and natural air pollution. In this context, Early Warning Systems (EWSs) play a key role in protecting the public during hazardous natural events, including DDSs.

EWSs play a critical role in protecting populations from climate-related and natural hazards, including DDSs. EWSs help reduce mortality and economic losses associated with extreme weather and hydrometeorological events. However, major implementation gaps persist, particularly in vulnerable regions such as small island developing states and least-developed countries. According to the United Nations, only 50% of countries operate adequate multi-hazard EWSs, despite the fact that 70% of disaster-related deaths over the past 50 years occurred in the most vulnerable areas. Bridging this gap by advancing forecasting technologies and ensuring broader access to EWSs remains essential, especially for managing transboundary phenomena like DDSs, where timely and accurate forecasts can significantly reduce health and environmental impacts.

The accuracy of EWSs ultimately depends on the performance of dust forecasting models, which vary considerably in their configuration and resolution. These models differ in first-layer calculation heights, ranging from approximately 20 to 100 meters above ground or sea level and in the particle size ranges they consider, spanning 0.03 to 20 μm (Knippertz and Stuut, 2014). They also incorporate diverse emission schemes (e.g., Marticorena and Bergametti, 1995) and deposition processes (e.g., Zender et al., 2003). These differences directly influence model outputs and contribute to substantial variability in dust concentration forecasts. Model evaluations often rely on sun photometers (e.g., AERONET; Holben et al., 1998) or satellite-derived aerosol optical depth (AOD) products (e.g., MODIS; NASA, 2024), which primarily capture atmospheric column dust loads rather than near-surface concentrations. These approaches, while valuable, suffer from key limitations. For example, cloud cover reduces satellite retrieval accuracy, and the low temporal resolution of polar-orbiting satellites limits event-scale analysis (Kazadzis et al., 2009). Despite these challenges, projects such as the Horizon Europe CiROCCO initiative aim to overcome these issues by integrating predictive frameworks and enhancing dust monitoring and forecasting capabilities.

Evaluating dust surface concentrations forecasts using ground PM_{10} concentrations is also challenging, as it is difficult to isolate the desert dust fraction in these measurements (Garcia-Castrillo & Tarradellas, 2017). Consequently, very few studies have evaluated numerical models using data from near-ground monitoring stations, and, to the best of our knowledge, none have focused on the Eastern Mediterranean Region (EMR). This region lies at the crossroads between Africa and Middle East, both of which are major sources of transboundary dust pollution. The study of the EMR has become increasingly difficult in the past decade due to ongoing conflicts and socio-political instability in the nearby dust source regions. This instability has limited in-situ data collection and delayed any mitigation efforts for DDSs (Eleftheriou et al., 2023).

In this study, we evaluate surface dust concentration forecasts from eleven (11) operational dust models and their multi-model ensemble (Multi-Model Median; MMM) using daily ground-level PM measurements from three background stations in the EMR. These stations, Ayia Marina (AM; Cyprus), Finokalia (FKL; Greece), and Be'er Sheva (BS; Israel), offer low-background



environments, allowing for reliable assessment of long-range transboundary dust transport. Observed PM data are categorized into PM₁₀, PM_{2.5}, and coarse particle fractions, and we apply a suite of statistical metrics to assess model performance. The measurements include 24-hour averages of both observed PM and modelled dust surface concentrations, ensuring temporal alignment. The performance evaluation incorporates different subsets of the data, including the entire study period, the 95th percentile of PM concentrations, and dust storm days, as classified using the methodology of Achilleos et al. (2020). Section 2 describes the EMR context and the selected study sites. Section 3 outlines data sources, model configurations, and the statistical evaluation methods. Section 4 presents results across all evaluation scenarios and PM fractions, including graphical diagnostics such as Taylor diagrams, BOOT methodology, and contingency tables. Section 5 summarizes findings and discusses implications for model refinement and operational forecasting.

The aim of this work is to determine how accurately current operational models forecast surface dust levels in the EMR, using background PM observations as reference. To our knowledge, this represents the first large-scale, multi-model evaluation of its kind for the region, which is critically exposed to dust intrusions but remains underrepresented in model validation studies. The findings highlight key model strengths and limitations and inform future efforts to improve dust prediction in support of health-focused EWS.

2. Characterizing DDSs in the EMR: Sources and Monitoring

The Mediterranean Basin frequently experiences DDSs, particularly in its southern and eastern regions, which are strongly influenced by emissions from Sahara and Sahel (Querol et al., 2009). Using air parcel back-trajectory models, Varga et al., (2014) demonstrate that dust often enters the Eastern Mediterranean Region (EMR) from Northern Africa, highlighting the dominant role of African sources in regional dust intrusions. Numerous studies report significant air quality degradation during these events in areas such as mainland Greece, the North Aegean, Cyprus, Israel, and Turkey (Çapraz et al., 2021; Triantafyllou et al., 2020; Tsiflikiotou et al., 2019; Vratolis et al., 2019; Krasnov et al., 2016; Mouzourides et al., 2015).

In addition to African sources, dust emissions from the Middle East significantly affect the EMR. For example, Bodenheimer et al. (2019) analysed 53 DDSs events, exceeding 150 µg/m³ using data from 13 air quality monitoring stations in Israel (2007–2013).

Their findings confirm the transboundary and multi-source nature of dust episodes, which often combine contributions from North African and Middle Eastern source regions.

Figure 1 illustrates the study area and background monitoring sites. These include Ayia Marina Xyliatou (AM) in Cyprus and Finokalia (FKL) in Crete, both part of the European Monitoring and Evaluation Programme (EMEP), and Be'er Sheva (BS) in Israel, part of the Israeli National Air Monitoring Network. These stations are strategically located in relatively isolated areas, minimizing local anthropogenic influence and making them ideal for evaluating long-range dust transport across the EMR.

Each station continuously records PM concentrations using high temporal resolution instruments: the Tapered Element Oscillating Microbalance (TEOM) for Cyprus and Israel, and the FH 62 I-R Thermo analyser for Greece. Measurements are collected at heights of 2–5 m above ground level at AM and FKL, and 10–15 m above ground at BS, where instruments are installed on building rooftops. These setups ensure consistent and representative sampling of ambient air, particularly for coarse PM fractions, which are most affected by dust intrusions.



Figure 1: Map of the study area showing the locations of the three background air quality monitoring stations: Ayia Marina Xyliatou (AM) in Cyprus, Finokalia (FKL) in Crete, Greece, and Be'er Sheva (BS) in Israel. The stations are strategically positioned away from major anthropogenic emission sources, making them suitable for evaluating transboundary desert dust transport. Imagery is obtained from © Google Earth.

In this study, we analyse daily averaged forecasts from 11 operational dust models available through the SDS-WAS platform, now known as the WMO Barcelona Dust Regional Center (<https://sds-was.aemet.es/>). Additionally, we use a Multi-model Median ensemble (MMM), which represents the median surface dust concentration predicted by all models for each site. Table 1 and

Table 2 summarise each model's key configurations like meteorological drivers, emission schemes, and spatial configuration, providing a basis for performance comparisons.

We compare modelled surface dust concentrations to daily observations of PM_{10} and $PM_{2.5}$ collected at the three sites. Before comparing with model outputs, the ground-based PM measurements underwent additional processing to ensure data consistency and reliability. This involves applying quality control checks to remove outliers and erroneous values, converting high-frequency data into 24-hour means, and harmonizing timestamps across sites to correct for time zone and logging differences. These steps are essential for minimizing bias and enabling a robust model–observation comparison.

The observation period covers 2012–2018 for AM, 2012–2017 for BS, and 2012–2016 for FKL. $PM_{2.5}$ data are not available for FKL during this timeframe. Data collection and processing follow the MEDEA project protocol (Achilleos et al., 2020; 2023), co-funded by the LIFE 2016 Programme, which aims to evaluate indoor exposure reduction interventions during DDSs. As part of this effort, Achilleos et al. (2020) develop a dust classification methodology that combines satellite imagery, $PM_{10}/PM_{2.5}$ ratios, elemental composition (Ca, Al, Fe), and dust aerosol optical depth (Dust-AOD) to standardize dust storm identification across all three sites. Using this multi-criteria approach, the authors identify 106 dust storm days at AM, 88 at FKL, and 101 at BS. Figure 2 presents the time series of daily PM_{10} concentrations for all three stations, with colour coding (black for AM, blue for FKL, red for BS) to distinguish each location. These records form the basis for assessing model performance across baseline, extreme, and dust-identified events, as detailed in subsequent sections.

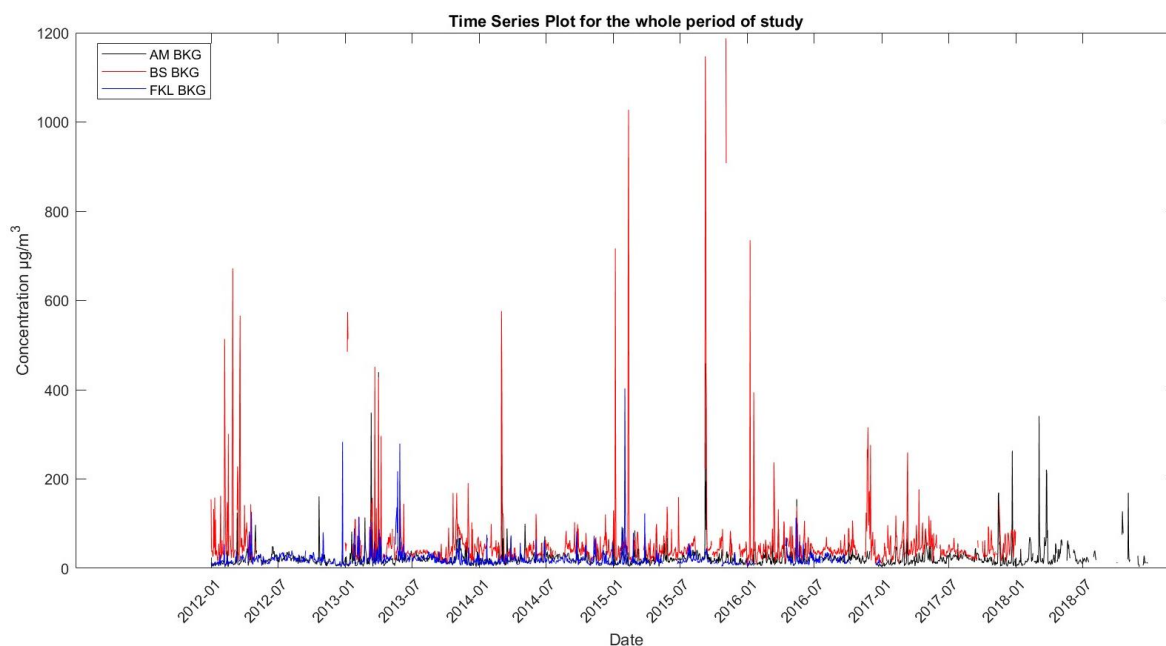


Figure 2: Time series of daily mean PM_{10} concentrations ($\mu\text{g}/\text{m}^3$) at the three monitoring stations over their respective study periods. Black represents Ayia Marina (Cyprus), blue represents Finokalia (Greece), and red represents Be'er Sheva (Israel). These measurements are used to evaluate dust model performance across baseline and extreme dust conditions.

140



Table 1: Configuration details of the 11 dust forecast models included in this study. These configurations correspond to the versions used during the study period (2012–2018) and may differ from current operational settings.

Model	Institution	Meteorological driver	Regional or global coverage	Meteorological initial fields	Radiation interactions	Emission scheme
BSC_DREAM8b_V2	BSC	Eta/NCEP	Regional	NCEP/GFS	Yes	Uplifting (Shao et al., 1993; Janjic et al., 1994; Ginoux et al., 2001)
MACC-ECMWF	CAMS	ECMWF	Global	ECMWF/IFS	No	Uplifting (Ginoux et al., 2001; Morcrette et al., 2008, 2009)
DREAM8-MACC	SEEVCCC	NMME	Regional	ECMWF/IFS	No	Uplifting (Shao et al., 1993; Janjic et al., 1994)
NMMB-BSC	BSC	NMMB/NC EP	Regional	NCEP/GFS	Yes	Saltation and sandblasting (Janjic et al., 1994; Marticorena and Bergametti, 1995)
NASA_GEOS	NASA	GEOS-5	Global	GEOS-5 Analysis	Direct effects fully included	Based on Ginoux (2001)
NCEP_NGAC	NCEP	NEMS GFS	Global	NCEP/GDA S	Yes (not active)	Dust uplifting following Ginoux (2001)
EMA_REG4	EMA	RegCM4	Regional	NCEP/GFS	Yes (both short and long waves)	Saltation and sandblasting (Zakey et al., 2006; Marticorena and Bergametti, 1995; Alfaro and Gomes, 2001)
DREAMABOL	CNR-ISAC	BOLAM	Regional	NCEP/GFS	Not active	Uplifting (Tegen et al., 1994)
NOA_WRF	NOA	WRF	Regional	GFS	No	GOCART scheme by Ginoux et al. (2001)
SILAM	FMI	ECMWF	Global	Offline Model	Photolysis rates are dependent on overlaying PM in our global 0.2 x 0.2 forecast with chemist	Sofiev et al., 2015
LOTOS	TNO	ECMWF	Regional	ECMWF	No	Marticorena & Bergametti 1995. Dust uplifting following Shao et al. (2001).



Table 2: Output configuration of the dust models used in the performance evaluation. These characteristics influence dust dispersion representation and forecast accuracy at ground level.

Model	Horizontal resolution	Vertical resolution	Height first layer	Transport size bins	Data assimilation	Reference
BSC_DREAM8b_V2	1/3° x 1/3°	24 Eta-layers	86 m (above sea level)	8 bins (0.1-10µm)	No	Nickovic et al., 2001; Pérez et al., 2006; Basart et al., 2012
MACC-ECMWF	8-10 km approx. (O1280)	137 sigma-layers	10 m (above surface)	3 bins (0.03-20µm)	Yes AOD550/MODIS	Rémy et al., 2022
DREAM8-MACC	1/3° x 1/3°	24 Eta-layers	96 m (above sea level)	8 bins (0.1-10µm)	Yes (ECMWF dust-analysis)	Nickovic et al., 2016
NMMB-BSC	1/3° x 1/3°	24 sigma-hybrid-layers	108 m (above surface)	8 bins (0.1-10µm)	No	Pérez et al., 2011; Klose et al., 2021
NASA_GEOS	0.25° x 0.3125°	72 layers (Top: 0.01 hPa)	-	5 bins (0.73-8.0 µm)	Yes (MODIS)	Colarco et al., 2010
NCEP_NGAC	T126 (~ 1°)	64 sigma-pressure hybrid layers (Top at 0.2 hPa)	20 m	5 bins (0.73-8.0 µm)	No	Lu et al., 2016 Zhang et al., 2022
EMA_REG4	45 km x 45 km	18 sigma-layers	50 m	4 bins (0.01-20)	No	Zakey et al., 2006
DREAMABOL	0.4° x 0.4°	50 sigma-hybrid levels	27 m above surface	8 bins (0.1-10µm)	No	Micrea et al., 2008; Mauruizi et al., 2011
NOA_WRF	0.19° x 0.22°	40 vertical levels	-	5 bins (0.73-8µm)	No	Flaounas et al., 2017
SILAM	0.1 x 0.1 deg (dust only) 0.2 x 0.2 deg (version with all modeled aerosols and full model chemistry)	19 hybrid levels up to about 16 km (0.1 x 0.1 dust only simulation) 28 hybrid levels up to about 50 km (0.2 x 0.2 simulation)	10m	4 bins (0.1-30 µm)	No	Sofiev et al., 2015
LOTOS	0.50° longitude × 0.25° latitude	12-15 layers	25 m above surface level	5 bins (0.1-10µm)	No	Manders et al., 2017

3. Evaluation Methodology and Performance Metrics

150 To evaluate the performance of the studied models, we use different statistical metrics, including the Pearson correlation coefficient (R), Mean Bias (MB), and Root Mean Square Error (RMSE). These metrics capture different aspects of model accuracy: R quantifies the strength of the linear relationship between observed and predicted values (ranging from -1 to 1), MB indicates



systematic over- or underestimation, and RMSE reflects the overall magnitude of error, with greater sensitivity to large discrepancies. The evaluation is carried out by comparing model outputs with observational data from the three background stations located in Cyprus, Greece, and Israel. This comprehensive assessment of dust models in the Eastern Mediterranean Region (EMR), to the best of our knowledge has never been done before and provides valuable insights into their performance in an understudied area.

3.1. Statistical Performance Metrics for Model Evaluation

The R coefficient measures the linear relationship between modelled and observed values, ranging from -1 to 1. A value of 1 indicates a perfect positive correlation, -1 indicates a perfect negative correlation, and a value of 0 implies no correlation (Carslaw, 2015). This metric is crucial for identifying whether the patterns of observation data and forecasts align.

$$R = \frac{1}{(n-1)} \sum_{i=1}^n \left(\frac{M_i - \bar{M}}{\sigma_M} \right) \left(\frac{O_i - \bar{O}}{\sigma_O} \right) \quad (1)$$

where M_i and O_i are the predicted and observed values respectively, and n is the number of observations.

Mean Bias (MB) is used to determine the average deviation of predicted values from observed values. A positive MB indicates that the model overestimates the observed values, while a negative MB indicates underestimation. This metric is significant as accurate forecasting can help vulnerable populations take precautionary measures during DDS events. The MB is defined as:

$$MB = \frac{1}{n} \sum_{i=1}^n M_i - O_i \quad (2)$$

The Root Mean Square Error ($RMSE$) measures the differences between predicted and observed values and is useful for indicating the model's overall prediction error. $RMSE$ penalizes larger errors more heavily than smaller ones, making it a valuable metric for assessing model performance. It is defined as:

$$RMSE = \left(\sum_{i=1}^n \frac{(M_i - O_i)^2}{n} \right)^{1/2} \quad (3)$$

While these metrics effectively summarize model performance, they also have limitations. Specifically, the R and $RMSE$ are linear measures and may fail to capture nonlinear relationships between observed and predicted values. Moreover, although $RMSE$ penalizes larger errors, it remains sensitive to extreme outliers, which potentially skewing the evaluation results, especially during intense dust events. Recognizing these limitations is crucial for interpreting the results and identifying areas where models may need further refinement.

3.2. Graphical Representations: Taylor Diagrams, BOOT Evaluation and Contingency Tables

Additionally, we employed the Taylor diagram, which provides a graphical summary of the three (3) metrics of R , $RMSE$ and the Standard Deviation (SD) in a single plot. Beyond summarizing performance metrics, the Taylor diagram offers a clear and intuitive means of comparing different models' performances against observational data. In this visualization the SD is proportional to the radial distance from the origin, while $RMSE$ corresponds to the distance from the “observed” point on the x-axis. The green dashed contours indicate specific values and the R is represented along the arc of the diagram, increasing from the top toward the bottom, as illustrated in Figure 3, for example.

The BOOT Model Evaluation Methodology (Chang & Hanna, 2004) is also used, incorporating bootstrapping techniques to resample the data and assess model performance in terms of underestimation or overestimation compared to observed values. The BOOT methodology was chosen over other resampling techniques, such as cross-validation, due to its flexibility and robustness in handling smaller datasets and its ability to provide a statistically sound assessment of model scatter and bias. Bootstrapping allows for the recreation of relationships between observed and modelled data using smaller, resampled datasets (Lahiri, 2010), providing a statistically robust comparison of model performance, and has been used in several recent studies (Lumet et al., 2024;



190 Yang et al., 2024). This methodology uses two key metrics: Geometric Mean (MG) a measure of relative bias and the Geometric Variance (VG) a measure of relative scatter, with the perfect model positioned at the bottom centre of the graph.

$$MG = \exp(\overline{\ln O_i} - \overline{\ln M_i}) \quad (4)$$

$$VG = \exp(\overline{\ln O_i - \ln M_i}^2) \quad (5)$$

195 Finally, categorical statistics were used to evaluate the effectiveness of models in capturing data points, rather than just assessing the intensity of values. A contingency table was developed to display how often models underestimated, overestimated, or accurately predicted PM concentrations (within 1-SD of the ground observations). The numbers in the heatmap represent the percentage of each condition (overestimation, underestimation, or accurate prediction) relative to the total number of model and observation comparisons for each day.

3.3. Approaches for Assessing Model Prediction against Observations

200 In this study, we evaluate model performance using multiple data subsets and comparative approaches to capture different aspects of forecast skill. The observed data include the PM_{10} ($O_{PM_{10}}$), $PM_{2.5}$ (O_{fine}) and PM_{coarse} (O_{coarse}) observations. The latter is calculated as the difference between PM_{10} and $PM_{2.5}$. For the model predictions, only dust concentration data ($M_{dustconc}$) is used. We apply two primary comparison strategies:

- i. Modelled dust concentration against PM_{10} observations ($M_{dustconc}$ vs $O_{PM_{10}}$),
- 205 ii. Modelled dust concentration against PM_{coarse} observations ($M_{dustconc}$ vs O_{coarse}).

These approaches enable evaluation of total particulate loading and the coarse-mode fraction, which is more directly linked to mineral dust. Each comparison is performed across three temporal evaluation scenarios:

- i. the entire study period for each region,
- ii. the 95th percentile of observed daily PM concentrations, representing extreme events, and
- 210 iii. identified dust storm days, based on the Achilleos et al. (2020), classification framework, which incorporates criteria such as PM_{10} exceedances, Dust Aerosol Optical Depth (Dust-AOD), MODIS satellite imagery, PM ratios, and elemental crustal signatures (Ca, Al, Fe).

This methodology follows the structure proposed by García-Castrillo & Terradellas (2017) but extends it by incorporating the BOOT evaluation technique (see Section 3.2). To our knowledge, this is the first multi-model validation effort
215 in the Eastern Mediterranean Region (EMR) to apply this combination of methods across multiple PM fractions, locations, and dust-intense periods. For clarity and completeness, we exclude outliers in BOOT visualizations that fall outside the evaluation graph boundaries, as these indicate very low predictive skill and would disproportionately affect interpretability.

4. Assessment of model performance

In this section, we present a comprehensive evaluation of model performance based on the statistical and graphical methods
220 described in Section 3. The results are organized around the three evaluation scenarios: (a) the entire study period, (b) the 95th percentile of observed PM concentrations, and (c) the subset of dust storm days identified using the Achilleos et al. (2020) criteria. Each model's ability to reproduce observed ground-level PM concentrations is assessed using PM_{10} and O_{coarse} data across all sites. Results are reported using statistical metrics (R , MB , and $RMSE$) and visual diagnostics, including Taylor diagrams, the BOOT methodology, and categorical heatmaps based on contingency analysis. This multi-metric framework allows for a nuanced
225 comparison of model skill, emphasizing both trend representation and bias magnitude under different dust-loading conditions.

4.1. Evaluation of Model Performance Under Different Scenarios



4.1.1. Results for the Entire Study Period

Table 3 presents a summary of model performance across the full study period, based on PM_{10} and O_{coarse} observations. Below, we present a comparative analysis of individual models, highlighting key trends and site-specific variability.

DREAM8-MACC achieves moderate correlations, with $R = 0.54$ at FKL, 0.51 at BS, and 0.48 at AM. It underestimates PM_{10} at AM ($MB = -13.87 \mu g/m^3$) and FKL ($MB = -13.55 \mu g/m^3$) but slightly overestimates at BS ($MB = 0.57 \mu g/m^3$). RMSE values are $26.05 \mu g/m^3$ (AM), $21.93 \mu g/m^3$ (FKL), and $79.26 \mu g/m^3$ (BS), indicating larger prediction errors in Israel. For O_{coarse} , MB shifts to $-3.43 \mu g/m^3$ (AM) and $21.84 \mu g/m^3$ (BS), suggesting sensitivity to particle fraction. The BOOT methodology is providing a good representation for the BS area with regards to the performance metrics whilst for AM it shows more underestimation than what the performance metrics indicate.

DREAM8b_V2 also shows moderate correlation ($R = 0.58$ at FKL, 0.49 at BS, 0.40 at AM), underestimating PM_{10} at AM ($MB = -7.93 \mu g/m^3$) and FKL ($MB = -10.32 \mu g/m^3$), and overestimating at BS ($MB = 0.57 \mu g/m^3$). RMSE ranges from $24.08 \mu g/m^3$ (FKL) to $85.20 \mu g/m^3$ (BS). For O_{coarse} , MB increases to $2.60 \mu g/m^3$ (AM) and $56.95 \mu g/m^3$ (BS), highlighting model sensitivity to coarse particle loading. The BOOT methodology graph aligns with the performance metrics, particularly for AM and BS.

DREAMABOL performs best at FKL ($R = 0.63$), with lower values at BS (0.49) and AM (0.42). It slightly underestimates PM_{10} at AM ($MB = -4.04 \mu g/m^3$) and FKL ($MB = -7.25 \mu g/m^3$), while overestimating at BS ($MB = 3.66 \mu g/m^3$). RMSE is highest at BS ($69.35 \mu g/m^3$). In the PM_{coarse} comparison, MB increases at AM ($6.32 \mu g/m^3$) and decreases slightly at BS ($23.31 \mu g/m^3$). The BOOT methodology for *DREAMABOL* shows consistency with the performance metrics described above except from the AM situation where even though MB changes from negative to positive between the approaches on the BOOT graph is still on the underestimation side.

EMA-REG4 shows the weakest performance, with $R = 0.31$ at FKL, 0.25 at BS, and just 0.02 at AM. It heavily underestimates at AM ($MB = -19.54 \mu g/m^3$) and FKL ($MB = -12.48 \mu g/m^3$), but overestimates at BS ($MB = 16.35 \mu g/m^3$). RMSE peaks at $172.40 \mu g/m^3$ (BS). Under O_{coarse} , MB improves at AM ($-8.77 \mu g/m^3$) but worsens at BS ($36.53 \mu g/m^3$), while RMSE remains high ($177.69 \mu g/m^3$). The poor performance of the model is also depicted in the BOOT methodology graph since the points for AM and FKL are out of the boundaries. Also, the points for BS are showing underestimation while the MB is positive.

LOTOS exhibits low correlation, with $R = 0.24$ – 0.47 , and consistent underestimation at all sites: $MB = -2.74 \mu g/m^3$ (AM), $-2.63 \mu g/m^3$ (BS), $-12.33 \mu g/m^3$ (FKL). RMSE values are highest at BS ($144.44 \mu g/m^3$) and AM ($46.39 \mu g/m^3$). Under O_{coarse} , MB reverses to $8.21 \mu g/m^3$ (AM) and $13.91 \mu g/m^3$ (BS), indicating a model bias toward larger particle overestimation. The low correlation of the model is also found on the BOOT methodology graph as the points for all sites are out of boundaries except from the BS point of the O_{coarse} approach that shows significant overestimation.

MACC-ECMWF performs moderately across all sites ($R = 0.47$ – 0.54), with negative MBs: $-11.84 \mu g/m^3$ (AM), $-28.50 \mu g/m^3$ (BS), $-11.51 \mu g/m^3$ (FKL). RMSE is highest at BS ($72.78 \mu g/m^3$). Under O_{coarse} , MB reduces significantly to $-1.39 \mu g/m^3$ (AM) and $-7.83 \mu g/m^3$ (BS), suggesting better alignment when fine PM is removed. The BOOT methodology captures the decrease in the negative values between the approaches for both Ayia Marina and Be'er Sheva.

The *Multi-Model Median (MMM)* achieves good correlation overall, with $R = 0.64$ (FKL), 0.58 (BS), and 0.56 (AM). It tends to underestimate at AM ($MB = -9.84 \mu g/m^3$) and FKL ($MB = -10.88 \mu g/m^3$) and slightly overestimate at BS ($MB = 13.74 \mu g/m^3$). RMSE values are $23.55 \mu g/m^3$ (AM), $20.43 \mu g/m^3$ (FKL), and $71.79 \mu g/m^3$ (BS). Under O_{coarse} , MB at AM improves to near-neutral ($0.66 \mu g/m^3$), while increasing at BS ($34.78 \mu g/m^3$). In the BOOT methodology, MMM supports the results in the performance metrics with Ayia Marina O_{coarse} and the Be'er Sheva PM_{10} approaches to be better performing.

NASA-GEOS performs best overall, with values of $R = 0.71$ (AM), $R = 0.65$ (BS), and $R = 0.64$ (FKL). These high correlations indicate that NASA-GEOS effectively captures the trends in PM_{10} concentrations at all locations, though the relatively high RMSE



values, particularly at Be'er Sheva ($89.48 \mu\text{g}/\text{m}^3$), highlight challenges in predicting extreme dust concentrations. The model's strong performance can be attributed to its configuration settings which include detailed representation of global dust dynamics (not regional) and effective calibration for long-range dust transport (i.e. Data assimilation). In terms of MB, NASA-GEOS shows a slight overestimation in Cyprus (MB = $3.60 \mu\text{g}/\text{m}^3$) and Israel (MB = $14.98 \mu\text{g}/\text{m}^3$), while it performs exceptionally well at Crete-Greece, where the bias is nearly zero (MB = $0.64 \mu\text{g}/\text{m}^3$). The model's RMSE values are relatively high, with the largest error seen at Israel (RMSE = $89.48 \mu\text{g}/\text{m}^3$), followed by Greece (RMSE = $41.24 \mu\text{g}/\text{m}^3$) and Cyprus (RMSE = $35.07 \mu\text{g}/\text{m}^3$), indicating some challenges in predicting extreme PM_{10} concentrations despite its overall strong trend capture. The models perform similarly in both approaches ($O_{\text{PM}_{10}}$ and O_{coarse}) but with decreases in R and increases in MB and RMSE for the O_{coarse} dataset for both areas (AM and BS). This difference indicates that the models underperform when trying to predict concentrations of larger particles. In the BOOT methodology analysis, NASA-GEOS performs best in both $M_{\text{dust conc}}$ vs $O_{\text{PM}_{10}}$ and $M_{\text{dust conc}}$ vs O_{coarse} approaches, further highlighting its strong predictive capabilities.

NCEP-NGAC exhibits site-dependent performance, with R = 0.62 at FKL, 0.52 at BS, and just 0.27 at AM. MB is positive across sites (3.40 – $7.96 \mu\text{g}/\text{m}^3$), while RMSE varies from $29.24 \mu\text{g}/\text{m}^3$ (FKL) to $68.76 \mu\text{g}/\text{m}^3$ (BS). O_{coarse} results show elevated MB: $18.34 \mu\text{g}/\text{m}^3$ (AM), $28.91 \mu\text{g}/\text{m}^3$ (BS). The Boot methodology graph presents an a very good performance on all 3 sites for the PM_{10} approach in comparison with the rest of the models. In terms of the O_{coarse} approach there is a shift towards the overestimation side which is consistent with the results of the performance metrics for Cyprus and Greece.

NMMB-BSC records moderate correlation, with R = 0.56 (BS), 0.43 (AM), and 0.42 (FKL). MB is consistently negative: $-15.03 \mu\text{g}/\text{m}^3$ (AM), $-25.04 \mu\text{g}/\text{m}^3$ (BS), and $-14.91 \mu\text{g}/\text{m}^3$ (FKL). RMSE is highest at BS ($86.19 \mu\text{g}/\text{m}^3$). In the O_{coarse} scenario, MB improves: $-4.46 \mu\text{g}/\text{m}^3$ (AM), $-3.68 \mu\text{g}/\text{m}^3$ (BS).

NOA-WRF performs well, particularly at BS (R = 0.62), FKL (0.58), and AM (0.56). MB shows overestimation: $7.41 \mu\text{g}/\text{m}^3$ (AM), $36.54 \mu\text{g}/\text{m}^3$ (BS), $1.67 \mu\text{g}/\text{m}^3$ (FKL). RMSE values reach $105.99 \mu\text{g}/\text{m}^3$ at BS. In O_{coarse} comparisons, MB increases to $17.54 \mu\text{g}/\text{m}^3$ (AM) and $54.64 \mu\text{g}/\text{m}^3$ (BS), consistent with the model's known bias in high dust load conditions. The BOOT methodology analysis shows that NOA-WRF consistently performs well for Be'er Sheva, particularly in the $M_{\text{dust conc}}$ vs $O_{\text{PM}_{10}}$ approach, further validating its robustness despite the tendency to overestimate. On the other hand, the BOOT methodology shows for both Ayia Marina and FKL sites to be always underestimating the observed values.

SILAM performs poorly overall, with R = 0.49 (FKL), 0.40 (BS), and 0.38 (AM). MB is negative at all sites: $-7.31 \mu\text{g}/\text{m}^3$ (AM), $-8.33 \mu\text{g}/\text{m}^3$ (BS), $-9.22 \mu\text{g}/\text{m}^3$ (FKL). RMSE is highest at BS ($95.29 \mu\text{g}/\text{m}^3$). Under O_{coarse} , MB becomes positive: $3.14 \mu\text{g}/\text{m}^3$ (AM), $8.73 \mu\text{g}/\text{m}^3$ (BS), suggesting the model underrepresents fine particles. The BOOT methodology is not supporting any of the findings for this model as no points are shown within the boundaries of the graph.

Summarizing the results from the site-specific analysis, Ayia Marina exhibits the highest overall correlation values across most models, with NASA-GEOS performing particularly well (R = 0.71). NOA-WRF and the MMM also effectively capture dust trends at this site, suggesting that key dust sources are well-represented. In contrast, EMA-REG4 and NCEP show substantial limitations, with correlations as low as R = 0.02 and R = 0.27, respectively, indicating poor simulation of dust transport for this region. At Be'er Sheva, model performance varies significantly, with systematic overestimation observed in several models (e.g., NOA-WRF MB = $36.54 \mu\text{g}/\text{m}^3$, NASA-GEOS MB = $14.98 \mu\text{g}/\text{m}^3$). This variability likely reflects the site's proximity to major Middle Eastern dust sources, which increases forecast complexity and uncertainty. At Finokalia, models generally show stronger agreement with observations, with NASA-GEOS and MMM achieving R = 0.64. The site's relative isolation from anthropogenic pollution makes it well-suited for evaluating long-range transboundary dust transport. However, models like LOTUS continue to underperform (R = 0.23), indicating persistent difficulties in simulating dust concentrations accurately at this location. Comparing the $O_{\text{PM}_{10}}$ and O_{coarse} approaches, results do not change substantially in terms of correlation. However, mean bias values fluctuate, with some



models shifting from negative to positive MB and vice versa, depending on particle size fraction. These variations highlight uncertainties linked to particle composition and model sensitivity to fine vs. coarse PM. Overall, the findings underscore the need for improved model tuning, especially with respect to local conditions and dominant particle size distributions in the Eastern Mediterranean Region.

Table 3: Statistical performance metrics (correlation coefficient R, mean bias MB, and root mean square error RMSE) for modelled dust concentrations compared to observed PM_{10} and O_{coarse} values at three background stations (Ayia Marina – AM, Be’er Sheva – BS, and Finokalia – FKL) over the full study period.

Model	Station	$O_{PM_{10}}$			O_{coarse}		
		R	MB ($\mu\text{g}/\text{m}^3$)	RMSE ($\mu\text{g}/\text{m}^3$)	R	MB ($\mu\text{g}/\text{m}^3$)	RMSE ($\mu\text{g}/\text{m}^3$)
DREAM8_MACC	AM	0,48	-13,87	26,05	0.47	-3.43	17.96
	BS	0,51	0,57	79,26	0.49	21.84	80.00
	FKL	0,54	-13,55	21,93	-	-	-
DREAM8b_V2	AM	0,40	-7,93	27,16	0.41	2.60	22.94
	BS	0,49	0,57	85,20	0.46	56.95	92.51
	FKL	0,58	-10,32	24,08	-	-	-
DREAMABOL	AM	0,42	-4,04	27,53	0.42	6.32	25.95
	BS	0,49	3,66	69,35	0.48	23.31	65.61
	FKL	0,63	-7,25	19,72	-	-	-
EMA_REG4	AM	0,02	-19,54	34,27	0.01	-8.77	24.79
	BS	0,25	16,35	172,40	0.24	36.53	177.69
	FKL	0,31	-12,48	21,54	-	-	-
LOTUS	AM	0,47	-2,74	46,39	0.45	8.21	47.30
	BS	0,24	-2,63	144,44	0.30	13.91	143.73
	FKL	0,23	-12,33	13,11	-	-	-
MACC_ECMWF	AM	0,47	-11,84	26,56	0.46	-1.39	20.83
	BS	0,54	-28,50	72,78	0.51	-7.83	60.59
	FKL	0,50	-11,51	28,08	-	-	-
MMM	AM	0,56	-9,84	23,55	0.56	0.66	18.46
	BS	0,58	13,74	71,79	0.56	34.78	76.47
	FKL	0,64	-10,88	20,43	-	-	-
NASA_GEOS	AM	0,71	3,60	35,07	0.7	14.23	39.44
	BS	0,65	14,98	89,48	0.6	36.11	100.09
	FKL	0,64	0,64	41,24	-	-	-
NCEP	AM	0,27	7,88	54,03	0.27	18.34	56.16
	BS	0,52	7,96	68,76	0.49	28.91	69.94
	FKL	0,62	3,40	29,24	-	-	-
NMMB	AM	0,43	-15,03	28,17	0.42	-4.46	20.78
	BS	0,56	-25,04	86,19	0.51	-3.68	85.56

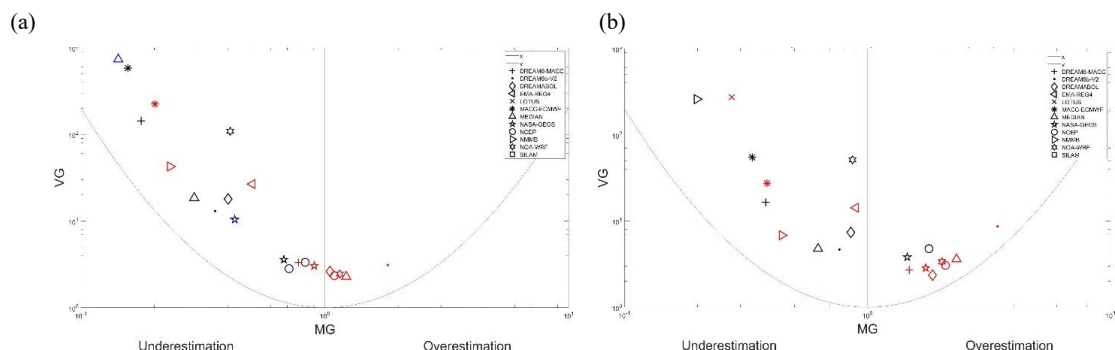


Figure 4: BOOT statistical plots comparing modelled daily PM_{10} concentrations to observed values over the full study period at the three monitoring stations: Ayia Marina (AM), Be'er Sheva (BS), and Finokalia (FKL). a) $M_{dust_{conc}}$ vs $O_{PM_{10}}$ and b) $M_{dust_{conc}}$ vs O_{coarse} . Models closer to the origin demonstrate better agreement with observations. Outliers beyond graph limits are excluded for visual clarity.

4.1.2. Results for the 95th Percentile of Measurements

This evaluation focuses on the models' ability to simulate extreme dust events, represented by the 95th percentile of daily mean PM concentrations. Results are presented for $O_{PM_{10}}$ and O_{coarse} , using correlation (R), mean bias (MB), and RMSE across the three monitoring stations.

NOA-WRF and NASA-GEOS rank as the top-performing models based on their consistently higher correlation coefficients and their ability to reproduce the temporal variability of extreme dust concentrations. Specifically, NOA-WRF achieves $R = 0.91$ at Finokalia, the highest correlation across all models and stations, reflecting its strong capacity to track peak dust events at this site. It also performs well at Ayia Marina ($R = 0.53$) and Be'er Sheva ($R = 0.62$), indicating robust spatial consistency. Similarly, NASA-GEOS shows strong correlations at Ayia Marina ($R = 0.66$) and Be'er Sheva ($R = 0.65$), and moderate performance at Finokalia ($R = 0.49$). These values suggest that NASA-GEOS more accurately represents the magnitude and timing of intense dust episodes at inland and island locations. Despite these strengths in correlation, both models exhibit notable positive bias, particularly at Be'er Sheva, where NOA-WRF and NASA-GEOS overestimate PM_{10} by $197.87 \mu g/m^3$ and $153.10 \mu g/m^3$, respectively, and show high RMSEs (302.59 and $300.83 \mu g/m^3$). These results indicate that while the models successfully capture dust event occurrence, they tend to exaggerate intensity, especially in proximity to dust source regions. The results from both models were supported by the BOOT methodology, which confirmed their relatively strong predictive capability, particularly in terms of capturing temporal trends (correlation), despite persistent overestimation during high dust events.

A group of moderate performers includes DREAMABOL, DREAM8-MACC, DREAM8b_V2, NCEP, and the MMM ensemble. These models exhibit moderate correlation values, typically ranging from $R = 0.20$ to 0.55 . For instance, DREAMABOL achieves $R = 0.60$ at FKL, while DREAM8b_V2 records $R = 0.45$ at BS. Biases vary: DREAMABOL and DREAM8b_V2 show underestimations at AM (MB = -24.00 and $-42.32 \mu g/m^3$, respectively), while NCEP shows mixed performance, with MB = $-86.47 \mu g/m^3$ at BS, but $+21.92 \mu g/m^3$ at FKL. The MMM model yields $R = 0.55$ at FKL and relatively balanced MB values (e.g., $-8.59 \mu g/m^3$ at AM, $17.00 \mu g/m^3$ at BS), but RMSE remains high across all locations (e.g., $252.92 \mu g/m^3$ at BS).

Models with consistently low performance were EMA-REG4, LOTUS, NMMB, and SILAM. EMA-REG4 showed $R = -0.24$ at AM and large underestimation (MB = $-90.34 \mu g/m^3$). At BS, despite $R = 0.41$, RMSE was $447.58 \mu g/m^3$, the highest among all models. LOTUS produced weak correlations across all sites, except in O_{coarse} at AM ($R = 0.59$). MB ranged from $-20.79 \mu g/m^3$ to $+98.75 \mu g/m^3$. NMMB systematically underestimated (MB = -50.52 to $-67.19 \mu g/m^3$) with poor correlation ($R < 0.4$) and RMSE up to $315.38 \mu g/m^3$. SILAM performed inconsistently, with $R = 0.10$ at AM, and although MB was near zero at FKL ($1.48 \mu g/m^3$),



RMSE exceeded $244 \mu\text{g}/\text{m}^3$. Among sites, Finokalia showed the best overall agreement, especially for NOA-WRF ($R = 0.91$), DREAMABOL (0.60), and MMM (0.55). Ayia Marina had high R for NASA-GEOS (0.66), but very poor scores for EMA-REG4 ($R = -0.24$). Be'er Sheva was the most challenging, with most models overestimating PM_{10} and RMSEs frequently exceeding $250 \mu\text{g}/\text{m}^3$. Using O_{coarse} instead of PM_{10} generally improved MB for several models. For instance, DREAM8-MACC at AM improved from $\text{MB} = -54.97$ to $-28.81 \mu\text{g}/\text{m}^3$, and NCEP at BS from -86.47 to $-22.56 \mu\text{g}/\text{m}^3$. However, RMSE reductions were limited, showing that errors in event intensity estimation persist even when fine particles are excluded.

In conclusion, while a few models (notably NOA-WRF and NASA-GEOS) demonstrated strong potential for predicting high-dust events, most models struggled with significant biases and poor precision. The coarse fraction approach slightly improved MB in many cases but did not resolve overall performance limitations. These results highlight the need for improved model representation of dust sources and transport, especially in complex environments like BS.

Table 4: Performance metrics (R , MB, RMSE) for each dust model during the 95th percentile of observed $O_{\text{PM}_{10}}$ and O_{coarse} concentrations at the three monitoring stations: Ayia Marina (AM), Be'er Sheva (BS), and Finokalia (FKL). The 95th percentile subset represents extreme dust events. MB and RMSE are reported in $\mu\text{g}/\text{m}^3$. Higher R values reflect better temporal agreement, while MB and RMSE quantify the magnitude and direction of forecast error. Results are shown separately for PM_{10} and $\text{PM}_{\text{coarse}}$.

Model	Station	$O_{\text{PM}_{10}}$			O_{coarse}		
		R	MB ($\mu\text{g}/\text{m}^3$)	RMSE ($\mu\text{g}/\text{m}^3$)	R	MB ($\mu\text{g}/\text{m}^3$)	RMSE ($\mu\text{g}/\text{m}^3$)
DREAM8_MACC	AM	0.20	-54.97	93.30	0.20	-28.81	68.70
	BS	0.40	-61.93	263.76	0.34	2.52	253.08
	FKL	0.53	-42.74	64.07	-	-	-
DREAM8b_V2	AM	0.22	-42.32	89.80	0.18	-12.69	70.91
	BS	0.45	-71.10	253.79	0.36	-9.57	239.84
	FKL	0.51	-14.46	63.21	-	-	-
DREAMABOL	AM	0.17	-24.00	79.28	0.10	4.25	69.74
	BS	0.39	-112.28	274.44	0.33	-50.58	226.17
	FKL	0.60	-15.90	54.98	-	-	-
EMA_REG4	AM	-0.24	-90.34	117.84	-0.23	-63.91	89.60
	BS	0.30	-53.59	447.58	0.24	0.10	444.59
	FKL	0.48	-46.97	55.23	-	-	-
LOTUS	AM	0.34	-20.79	108.52	0.59	-3.94	54.19
	BS	0.33	27.73	240.86	0.20	98.75	325.55
	FKL	-	-	-	-	-	-
MACC_ECMWF	AM	0.23	-43.39	90.59	0.20	-16.51	70.76
	BS	0.34	-157.67	286.51	0.28	-97.75	244.58
	FKL	0.40	-25.17	88.50	-	-	-
MMM	AM	0.32	-35.41	80.27	0.31	-8.59	60.33
	BS	0.42	-46.16	252.92	0.35	17.00	244.69
	FKL	0.55	-25.42	56.38	-	-	-



NASA_GEOS	AM	0.66	41.01	113.29	0.63	68.06	128.02
	BS	0.49	83.73	300.83	0.40	153.10	345.94
	FKL	0.56	44.96	127.04	-	-	-
NCEP	AM	0.32	-10.50	78.27	0.31	17.63	69.67
	BS	0.40	-86.47	251.04	0.30	-22.56	226.39
	FKL	0.52	21.92	67.91	-	-	-
NMMB	AM	0.18	-53.53	96.58	0.15	-25.29	75.86
	BS	0.44	-67.19	315.37	0.35	-7.00	324.01
	FKL	0.27	-50.52	81.93	-	-	-
NOA_WRF	AM	0.53	25.96	78.54	0.51	48.40	87.22
	BS	0.55	163.83	282.49	0.60	197.87	302.59
	FKL	0.91	37.45	66.98	-	-	-
SILAM	AM	0.10	25.01	244.42	0.05	56.54	245.77
	BS	0.39	11.73	228.50	0.33	43.18	234.93
	FKL	0.25	1.48	76.28	-	-	-

365

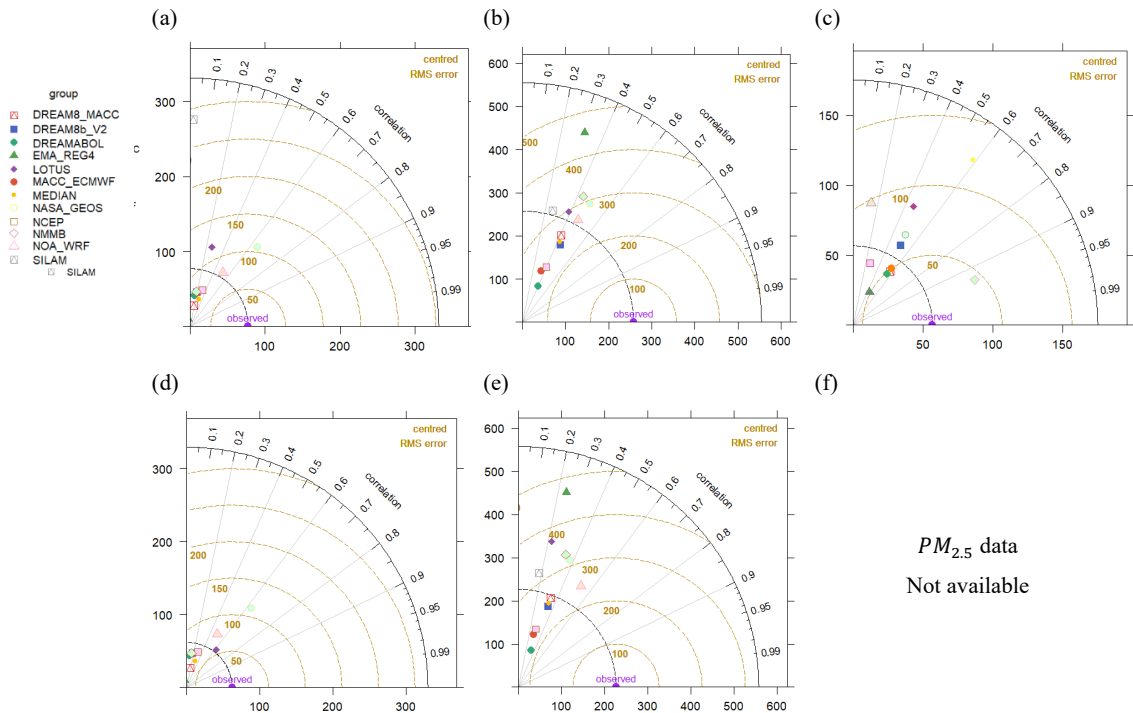


Figure 5: Taylor diagrams showing the performance of 11 operational dust models and the multi-model median (MMM) during the 95th percentile of observed PM_{10} concentrations at the three monitoring stations: (a and c) Ayia Marina (AM), (b and e) Be'er Sheva (BS) and (c and f) Finokalia (FKL). This analysis reflects the ability of each model to reproduce the timing and magnitude of extreme dust events.

370

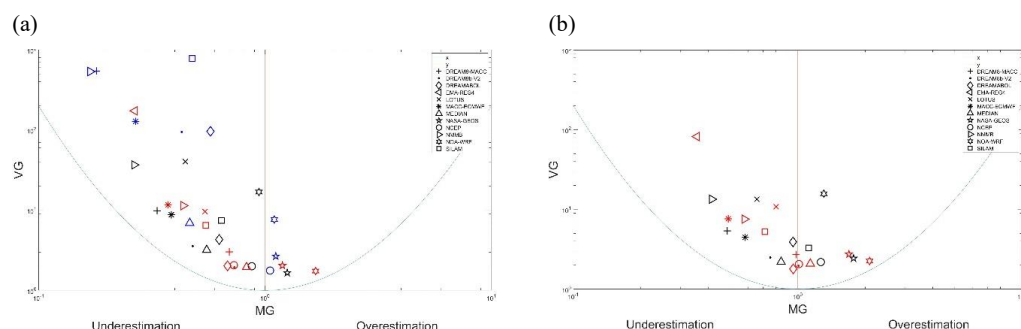


Figure 6: BOOT statistical plots comparing modelled and observed PM_{10} concentrations for the 95th percentile of daily values at Ayia Marina (AM), Be'er Sheva (BS), and Finokalia (FKL). a) $M_{dust_{conc}}$ vs $O_{PM_{10}}$ and b) $M_{dust_{conc}}$ vs O_{coarse} . Models positioned closer to the origin show better performance. Extreme outliers are excluded for readability.

4.1.3. Results for Dust Days Identified by Achilleos et al. (2020) methodology

The evaluation of dust days identified by Achilleos et al. (2020) highlights significant variability in model performance across the three sites—Ayia Marina (AM), Be'er Sheva (BS), and Finokalia (FKL). The results in this section will not emphasize changes in statistical values with deviations of less than 10%, as these are not considered significant in the overall evaluation of model performance and are excluded for brevity (Table 5).

At Ayia Marina, correlations are generally low, with only NASA-GEOS ($R=0.65$) and NOA-WRF ($R=0.51$) performing well. Most models underestimate dust levels, as reflected in negative mean bias (MB) values, except for NOA-WRF and SILAM, which show overestimations. The RMSE values remain high ($> 50 \mu g/m^3$), indicating substantial uncertainties during dust events. EMA-REG4 performs particularly poorly, exhibiting a negative correlation ($R=-0.07$) for this location.

At Be'er Sheva, model performance is mixed, with only four models—NOA-WRF, MMM, NASA-GEOS, and DREAM8b_V2—achieving correlations above $R=0.5$. Models like NASA-GEOS, SILAM, and NOA-WRF tend to overestimate PM concentrations, while most others exhibit underestimation. High RMSE values, often exceeding $150 \mu g/m^3$, reflect the site's proximity to major Middle Eastern dust sources, complicating accurate dust event representation.

At Finokalia, model correlations are relatively higher, with NOA-WRF ($R=0.68$) and DREAMABOL ($R=0.62$) performing best. The station's remote location, less influenced by local anthropogenic activities, allows models to better capture transboundary dust transport. However, overestimations remain common, particularly for NASA-GEOS ($MB=43.66 \mu g/m^3$).

Comparison of the two approaches ($O_{PM_{10}}$ vs. O_{coarse}) reveals minimal changes in correlations, but MB often shifts from underestimation to overestimation, particularly at Be'er Sheva and Ayia Marina. This suggests that finer particulate matter in PM_{10} is not adequately captured by models, leading to higher underestimations. The RMSE values remain high for all sites, underscoring the models' uncertainties during extreme dust events.

The BOOT methodology generally aligns with the MB and RMSE trends but shows notable inconsistencies for some models, particularly those with poor performance. For example, models such as EMA-REG4 and LOTUS display mismatches between their BOOT results and statistical metrics, with BOOT sometimes placing them on the underestimation side even when MB indicates overestimation. Additionally, the BOOT methodology does not fully capture the unique challenges of sites like Be'er Sheva, where proximity to major dust sources introduces higher RMSE values. For example, models such as NASA-GEOS and NOA-WRF perform well in terms of MB and R values but are depicted as having larger scatter on the BOOT graph due to the



higher RMSE values typical of this site. Despite these inconsistencies, BOOT provides valuable insights into overall model scatter and bias trends, highlighting areas for further methodological refinement.

Table 5: Statistical performance metrics (correlation coefficient R, mean bias MB, and Root Mean Square Error RMSE) for modelled dust concentrations compared to observed $O_{PM_{10}}$ and O_{coarse} values at three background stations (Ayia Marina – AM, Be’er Sheva – BS, and Finokalia – FKL) over for the Achilleos et al. (2019) approach.

Model	Station	$O_{PM_{10}}$			O_{coarse}		
		R	MB ($\mu\text{g}/\text{m}^3$)	RMSE ($\mu\text{g}/\text{m}^3$)	R	MB ($\mu\text{g}/\text{m}^3$)	RMSE ($\mu\text{g}/\text{m}^3$)
DREAM8_MACC	AM	0.22	-29.31	62.12	0.23	-9.77	44.95
	BS	0.47	-13.81	178.55	0.43	26.37	172.53
	FKL	0.48	-27.73	51.53	-	-	-
DREAM8b_V2	AM	0.22	-17.55	60.51	0.22	1.65	50.04
	BS	0.51	-1.05	165.96	0.44	39.91	162.37
	FKL	0.54	-7.83	51.72	-	-	-
DREAMABOL	AM	0.17	-4.15	56.11	0.18	14.04	51.11
	BS	0.41	-22.64	186.52	0.37	14.82	161.09
	FKL	0.62	-7.56	44.94	-	-	-
EMA_REG4	AM	-0.07	-48.03	69.43	-0.07	-29.63	50.50
	BS	0.28	5.96	325.89	0.25	43.98	333.95
	FKL	0.24	-29.62	50.77	-	-	-
LOTUS	AM	0.15	-7.20	81.04	0.18	11.13	77.39
	BS	0.08	27.86	258.60	0.10	50.90	253.05
	FKL	-	-	-	-	-	-
MACC_ECMWF	AM	0.26	-20.19	60.44	0.26	-0.43	49.57
	BS	0.42	-75.36	186.89	0.37	-35.71	159.61
	FKL	0.43	-17.38	69.81	-	-	-
MMM	AM	0.32	-14.36	53.71	0.34	5.29	42.89
	BS	0.51	7.10	167.23	0.47	47.51	166.29
	FKL	0.56	-12.89	44.88	-	-	-
NASA_GEOS	AM	0.65	24.70	79.13	0.64	44.48	89.58
	BS	0.59	75.92	212.60	0.52	118.21	241.70
	FKL	0.56	43.66	106.02	-	-	-
NCEP	AM	0.25	11.39	64.43	0.26	31.01	64.19
	BS	0.46	-10.07	167.43	0.39	29.14	157.88
	FKL	0.46	31.63	70.06	-	-	-
NMMB	AM	0.19	-28.27	65.52	0.20	-8.72	52.03
	BS	0.50	-37.31	211.68	0.43	2.08	216.35
	FKL	0.28	-33.26	63.45	-	-	-



NOA_WRF	AM	0.51	42.39	79.25	0.52	60.03	89.31
	BS	0.61	128.66	215.47	0.64	158.60	237.19
	FKL	0.68	40.68	64.90	-	-	-
SILAM	AM	0.20	28.18	191.44	0.19	45.78	195.98
	BS	0.26	32.42	189.39	0.23	58.12	198.82
	FKL	0.19	16.30	72.42	-	-	-

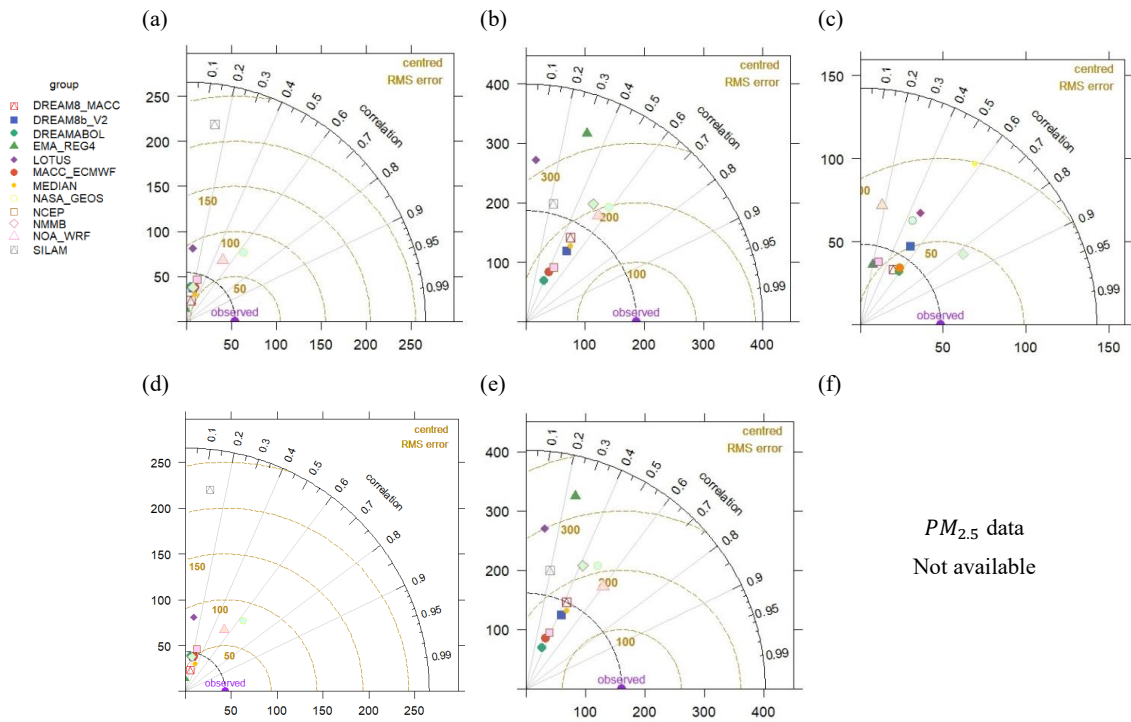


Figure 7: Taylor diagrams showing the performance of 11 operational dust models and the multi-model median (MMM) during identified dust storm days at the three monitoring stations: (a and c) Ayia Marina (AM), (b and e) Be'er Sheva (BS) and (c and f) Finokalia (FKL).

410

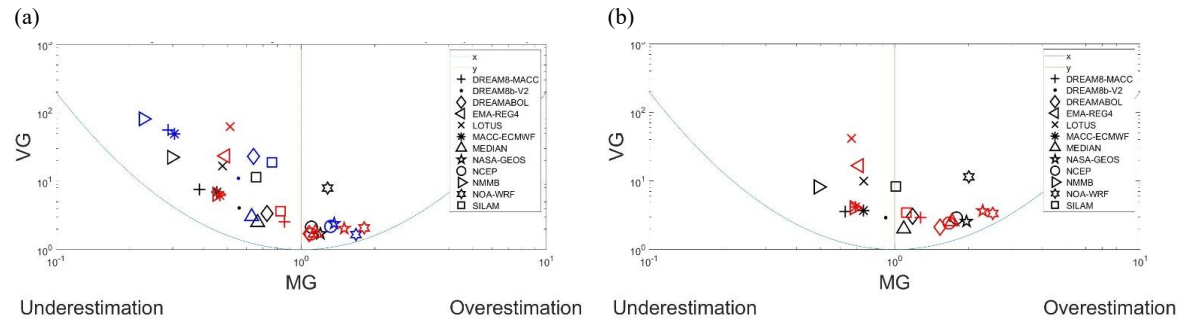




Figure 8: BOOT statistical plots comparing modelled and observed PM_{10} concentrations for the Achilleos et al. (2019) approach at Ayia Marina (AM), Be'er Sheva (BS), and Finokalia (FKL). a) $M_{dustconc}$ vs $O_{PM_{10}}$ and b) $M_{dustconc}$ vs O_{coarse} . Models positioned closer to the origin show better performance. Extreme outliers are excluded for readability.

Summarizing the results based on the Achilleos et al. (2019) approach and site-specific analysis, Ayia Marina shows low overall correlation values for most models, with only NASA-GEOS ($R = 0.66$) and NOA-WRF ($R = 0.53$) performing well. The remaining models fail to capture dominant dust patterns, indicating insufficient representation of key dust sources for this site. EMA-REG4 displays negative correlation ($R = -0.24$), clearly failing to simulate observed conditions. Most models tend to underestimate major dust events, with only a few (e.g., NASA-GEOS, SILAM) producing slight overestimations. In two of the three sites, models with higher correlation also exhibit stronger overall performance.

At Be'er Sheva, model performance is generally poor. Only NOA-WRF ($R = 0.62$), MMM ($R = 0.58$), NASA-GEOS ($R = 0.65$), and DREAM8b_V2 ($R = 0.45$) exceed $R = 0.5$. Underestimation dominates, though models like NASA-GEOS, SILAM, and NOA-WRF show positive mean bias values (e.g., MB = $197.87 \mu\text{g}/\text{m}^3$ for NOA-WRF). The site's proximity to Middle Eastern dust sources likely contributes to this variability, as reflected in the significantly higher RMSE values compared to Ayia Marina and Finokalia. At Finokalia, most models perform comparatively better. NOA-WRF ($R = 0.91$) and DREAMABOL ($R = 0.60$) achieve the highest correlations. The site's limited anthropogenic influence and exposure to transboundary dust enhance model agreement with observations.

Comparing the two evaluation approaches, $O_{PM_{10}}$ and O_{coarse} , reveals minor differences in R coefficient, but notable changes in MB. In several cases, MB shifts from negative to positive or vice versa. A general trend is observed: underestimating models reduce their negative bias, and slightly overestimating models exhibit larger positive bias under the O_{coarse} approach. This suggests that finer particulate matter in $O_{PM_{10}}$ is not fully captured, leading to stronger underestimations, particularly for models that do not resolve $PM_{2.5}$ well.

RMSE remains consistently high across all sites, confirming persistent model uncertainty during intense dust events. The BOOT evaluation supports the MB trends but shows discrepancies for low-performing models, further emphasizing the limitations in model reliability under extreme conditions.

4.2. Categorical statistics

The following results present the contingency tables in the form of heatmaps, depicting the daily agreement of modelled and observed values at each site. The results are categorized into three scenarios: overestimation, underestimation and accurate prediction (hit; within 1-SD of the ground observations) for each approach and dataset used. The two approaches include O_{coarse} and $O_{PM_{10}}$, focusing on how well the models capture different PM fractions at each site.

4.2.1. Ayia Marina (AM), Cyprus

Figure 9a and b show the heatmaps for the Ayia Marina site, comparing model performance for O_{coarse} and $O_{PM_{10}}$ over the entire study period. All models achieve over 80% hits for both approaches, demonstrating their ability to capture the overall dust conditions at this site. The DREAM8-MACC model stands out with the highest hit rate at 97%, indicating excellent performance in predicting both coarse and total PM concentrations. On the other hand, NOA-WRF records the lowest hit rate (81%), reflecting a tendency to overestimate the dust levels in this region, as confirmed by the model's consistent overestimation across multiple metrics in prior analyses. Underestimation remains minimal for Ayia Marina in both approaches, indicating that most models are capturing or over-predicting dust events rather than missing them. This is critical for public health warnings, as underestimation would lead to insufficient precautionary measures.



For the 95th percentile of ground observations (Figure 10a and b), AM shows mixed results between hits and overestimations in the O_{coarse} approach. NOA-WRF and NASA-GEOS display the highest overestimations, consistent with their tendency to over-predict extreme dust events, which has been observed across multiple metrics. EMA-REG4, on the other hand, records the fewest overestimations and is more prone to underestimating, continuing its underperformance noted in earlier sections. The overall hit percentages remain high for most models, but overestimations dominate during peak dust events. In the $O_{PM_{10}}$ approach, the results are similar, with NOA-WRF and NASA-GEOS again leading in overestimations, but no strong trends emerge beyond these individual cases. EMA-REG4 again shows high underestimation, making it the weakest performer during these peak dust events. For the specific dates examined, following the third approach of Subsection 3.3, the O_{coarse} approach shows (Figure 10a) that most models perform well, capturing ground concentrations within one standard deviation. EMA-REG4 has the highest hit percentage, which is notable given its underperformance in other scenarios. DREAM8-MACC also performs well, while NOA-WRF shows the highest overestimation rate, continuing the trend seen in previous sections. Underestimations remain low, indicating that most models effectively predict peak dust levels.

4.2.2. Be'er Sheva (BS), Israel

The performance at Be'er Sheva, shown in Figure 9a and b, is more varied compared to AM. For the O_{coarse} approach, the models exhibit a broader range of results. MACC-ECMWF and NMMB models show the highest hit percentages, indicating strong predictive capabilities for coarse particles. In contrast, DREAM8b-V2 has the lowest hit rate, suggesting a struggle to accurately capture dust levels at this location, particularly for coarse particles. In the $O_{PM_{10}}$ approach, the performance remains largely the same, with MACC-ECMWF and NMMB continuing to outperform other models. However, DREAM8b-V2 again records the lowest hit rate, reinforcing its overall weaker performance at this site, which was also observed in the earlier analyses of correlation and bias metrics. Overestimation is more frequent than underestimation, particularly for NOA-WRF, which tends to significantly overestimate PM concentrations, leading to potential overreactions in dust event management.

For the 95th percentile of ground observations, overestimations dominate the results for both approaches (O_{coarse} and $O_{PM_{10}}$), as depicted in Figure 10a and b. This site's proximity to major dust sources in the Middle East may contribute to the over-prediction, as models struggle to account for rapid dust inflow. NOA-WRF shows the highest overestimation, with no underestimations recorded for the O_{coarse} approach, reinforcing the model's consistent bias towards overestimating dust levels at this site. EMA-REG4, which previously showed underestimations at other sites, also demonstrates high underestimations for $O_{PM_{10}}$, indicating that it struggles to capture the dust dynamics in Be'er Sheva.

For the specific dates examined, at Be'er Sheva, overestimations dominate again in the O_{coarse} approach, with NOA-WRF and NASA-GEOS showing the highest overestimation percentages. EMA-REG4, typically prone to underestimations, continues to show a relatively balanced performance between hits and underestimations. In the $O_{PM_{10}}$ approach (Figure 11b), the results are consistent, with high overestimation rates, reinforcing the models' struggle to accurately predict dust levels in this challenging region.

4.2.3. Finokalia (FKL), Greece

At Finokalia, the models generally perform well, with all showing hit percentages above 88%, as illustrated in Figure 9. The LOTUS model performs exceptionally well in both approaches, achieving a 100% hit rate for O_{coarse} , indicating that it accurately captures dust conditions at this site. This high performance contrasts with LOTUS's weaker results at other sites, like Be'er Sheva, highlighting its region-specific strengths. Conversely, NCEP records the lowest hit percentage, suggesting challenges in predicting both coarse and total PM at Finokalia. Finokalia, with its more consistent results, supports the idea that this site, being less



influenced by anthropogenic sources, allows models to focus more on transboundary dust transport, leading to higher accuracy across most models. This observation aligns with the results from previous sections, where high correlations and lower biases were recorded for many models at this site.

490 For this area, there are no trends using the 95th percentile of ground observations, in either approach, with most models showing a balance between hits, overestimations, and underestimations (Figure 10). EMA-REG4 and NMMB record the highest underestimations, while NOA-WRF shows the highest overestimations, consistent with its overall behaviour across the sites. The hit percentages remain relatively high, but the variability indicates that extreme dust events in Finokalia are more challenging for models to predict accurately, likely due to the site's geographical characteristics.

495 Finally, for specific dates approach, and using $O_{PM_{10}}$ data, results show (Figure 11b) that most models achieve more than 50% hits, with DREAMABOL achieving the highest hit percentage. NOA-WRF remains the only model to exceed 80% overestimation, reinforcing its tendency to over-predict dust levels at this site. These results align with earlier observations, suggesting that while most models perform well at Finokalia, there is still room for improvement in accurately capturing peak dust events.

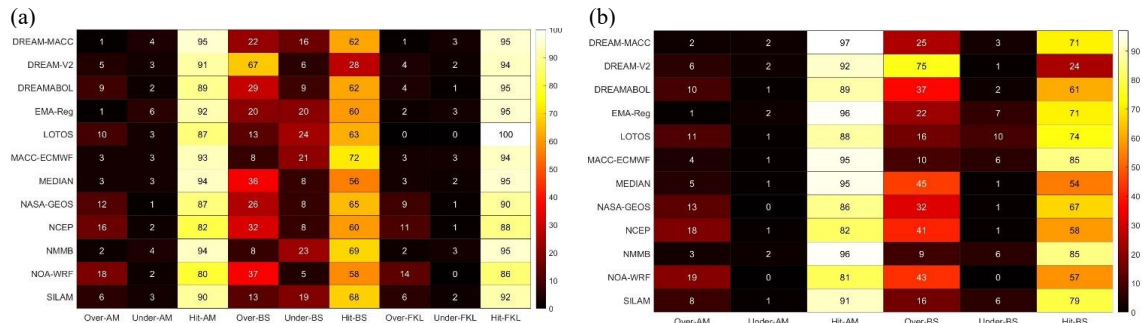
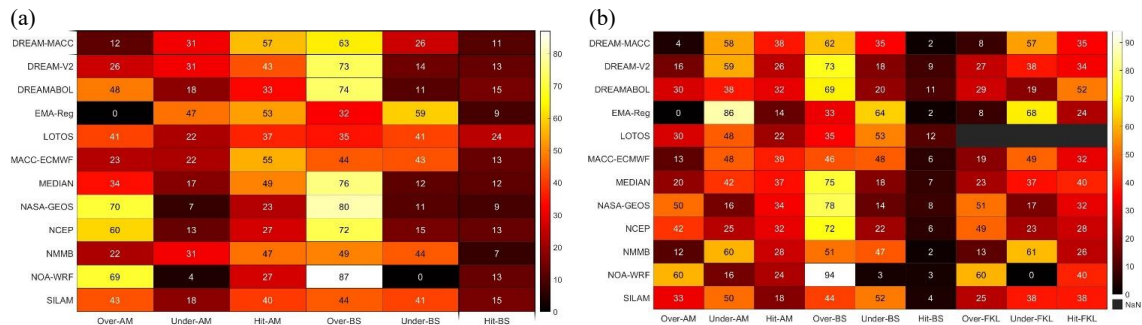


Figure 9: Heatmap of the scenarios (a) $M_{dust_{conc}}$ vs O_{coarse} and (b) $M_{dust_{conc}}$ vs $O_{PM_{10}}$ for the entire period of study



500 Figure 10: Heatmap of the scenarios (a) $M_{dust_{conc}}$ vs O_{coarse} and (b) $M_{dust_{conc}}$ vs $O_{PM_{10}}$ for the 95th percentile of data

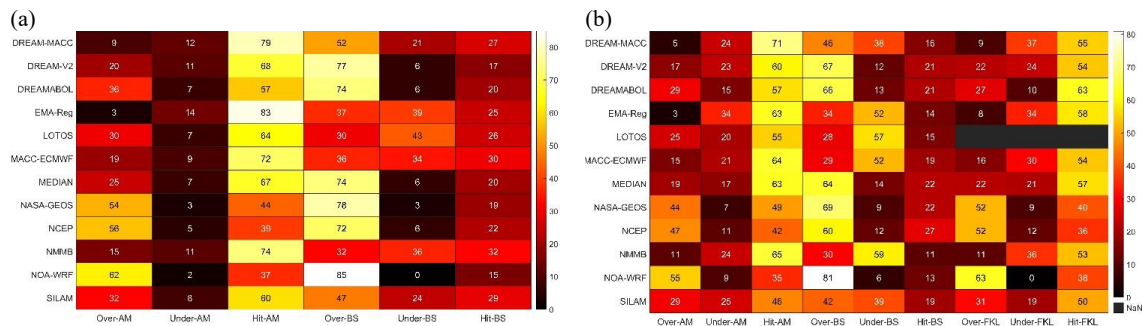


Figure 11: Heatmap of the scenarios (a) $M_{dust_{conc}}$ vs O_{coarse} and (b) $M_{dust_{conc}}$ vs $O_{PM_{10}}$ for specific dust days.



Overall, no model demonstrates exceptional accuracy across all sites, approaches and datasets. A high hit rate of models accurately predicting observations is observed only in the first dataset, which assesses performance over the entire study period. Interestingly, models like EMA-REG4 and NMMB, which consistently perform poorly in other evaluation methods, show significantly better results when evaluated using contingency table methods.

5. Evaluation of Configuration Settings in the performance of Operational Dust Forecasting Models

The performance of dust forecasting models is critically influenced by several key configuration parameters, as outlined in Table 1 (Model Input Configuration) and

Table 2 (Model Output Configuration). These parameters include vertical resolution, the height of the first layer, data assimilation techniques, horizontal resolution, meteorological drivers, and emission schemes. Each of these factors contributes to the models' ability to accurately predict dust events, which is essential for effective environmental and health risk management.

Models with higher vertical resolution, such as NASA-GEOS and MACC-ECMWF, consistently perform well across various study sites, largely due to their capability to accurately represent the vertical distribution of dust. This feature is especially crucial during long-range transport events, where models need to capture dust movement at different atmospheric levels. For instance, NASA-GEOS achieved strong correlation coefficients at sites like Ayia Marina and Finokalia, underscoring the importance of vertical resolution in complex dust scenarios. In contrast, models with fewer vertical layers, like LOTUS, show limitations in areas with complex topography or high variability in dust inflow. This variability highlights how model design affects performance depending on site characteristics.

Another key factor is the proximity of the first model layer to the ground, which has a direct impact on the accuracy of near-surface dust concentration predictions. Models such as MACC-ECMWF and NCEP-NGAC, with their first layers positioned closer to the surface, show higher correlations with observed particulate matter at surface level. This finding suggests that a lower starting height enhances a model's ability to capture near-surface dust levels accurately—particularly valuable during intense dust events when ground-level concentrations are of greatest concern.

Beyond structural design, data assimilation significantly enhances model performance, especially in adjusting predictions for real-time accuracy. Models incorporating data assimilation, like NASA-GEOS, achieve not only higher correlation coefficients but also reduced biases and RMSEs. This is especially beneficial during extreme dust events, where rapid updates based on real-time data are critical for maintaining predictive accuracy. The effectiveness of data assimilation across models further demonstrates its value in capturing the dynamic nature of dust events.

Horizontal resolution also plays an important role, as finer grids allow models to improve spatial accuracy, particularly in regions with complex terrain. For instance, models with finer grids, such as NOA-WRF, demonstrate improved ability to handle spatial variability, yielding higher correlation and lower bias at specific sites like Be'er Sheva and Finokalia. This reinforces the need for spatial detail when modelling dust in diverse geographical settings, where elevation and landscape variation influence dust distribution patterns.

Finally, the choice of meteorological driver and the specific emission schemes employed by models also significantly impact their performance. For example, models driven by ECMWF and GEOS-5 generally outperformed those with other meteorological drivers, demonstrating both higher correlation coefficients and lower biases. Similarly, emission schemes that effectively simulate dust lifting and deposition processes, such as those used in MACC-ECMWF and NASA-GEOS, tend to enhance model performance, particularly in accurately forecasting the intensity and spatial distribution of dust during high-dust events.



In conclusion, the configuration of dust models plays a crucial role in their predictive accuracy, with vertical resolution, the height of the first layer, data assimilation, horizontal resolution, meteorological drivers, and emission schemes all influencing their effectiveness. Understanding these relationships is crucial for improving model forecasts and effectively managing the environmental and health impacts of dust events in the Eastern Mediterranean region.

6. Conclusions and Future Directions in Dust Modelling in EMR

This study provides a comprehensive evaluation of eleven operational dust models and a multi-model ensemble (MMM) in forecasting PM concentrations across the EMR. Evaluating operational dust forecasting numerical models against ground measurements is essential to identifying for determining which models perform better in different areas. Such evaluations are crucial for enhancing the accuracy and reliability of DDS - EWS, as these systems rely on robust forecasting models to provide timely alerts for hazardous dust events. Despite their importance, few studies have conducted such evaluations against ground-level measurements, highlighting a critical gap that this research aims to address.

6.1. Key findings

The results clearly indicate that certain models outperform others in predicting ground-level PM concentrations across specific regions, although their performance varies based on the evaluation method applied. Notably, no single model consistently achieves accurate predictions across all three regions, underscoring the need for model-specific adaptations or improvements tailored to regional characteristics. The key findings of this study are summarised as follows:

1. Model Performance Across Sites: The accuracy of individual models varied significantly across the three sites (Ayia Marina, Finokalia, and Be'er Sheva). NASA-GEOS emerged as the most consistent performer, particularly at Ayia Marina and Be'er Sheva, benefiting from a high-resolution vertical structure, data assimilation, and robust meteorological inputs via the GEOS-5 driver. Conversely, models like EMA-REG4 showed substantial limitations, especially in accurately predicting PM levels in Cyprus, indicating that regional drivers and lower-resolution configurations may not be suitable for cross-boundary events in the EMR.

2. Enhancing Early Warning Systems: The findings of this study can provide valuable insights for improving EWS and mitigation strategies in the EMR. By identifying the strengths and weaknesses of the evaluated dust models, particularly in predicting extreme events and PM concentrations, these results can help tailor EWS to the region's specific needs. For example, models such as NASA-GEOS and NOAA-WRF, which demonstrated higher accuracy and correlation during intense dust events, can be prioritized for operational forecasting to provide timely warnings to vulnerable populations. Additionally, the multi-model ensemble (MMM), with its reduced uncertainty, can enhance the reliability of forecasts, supporting decision-makers in public health protection and emergency preparedness. These improved predictions are essential for issuing health advisories, minimizing exposure to hazardous PM levels, and informing policy measures to mitigate the adverse impacts of DDS events on air quality and public health.

3. Impact of Model Configurations on Performance:

i. **Meteorological Drivers and Data Assimilation:** ECMWF and GEOS-5-driven models demonstrated higher correlations with observed data, particularly when data assimilation was included. Data assimilation improves models' responsiveness to rapid dust influx and significantly enhances reliability in high-dust scenarios, underscoring its necessity in complex regions like the EMR. Optimizing dust forecasting models is crucial for improving predictive accuracy, especially in regions affected by transboundary dust transport. These advancements align with the goals of initiatives like the Horizon Europe CiROCCO project, which aims to strengthen dust storm monitoring through integrated predictive frameworks.



ii. **Horizontal and Vertical Resolution:** Models with finer horizontal and vertical resolutions, such as NOAA-WRF

($0.19^\circ \times 0.22^\circ$, 30 layers) and NASA-GEOS ($0.25^\circ \times 0.3125^\circ$, 72 layers), generally showed improved spatial and temporal accuracy across sites. Enhanced vertical resolution, particularly for capturing dust plumes at altitude, contributed to better event tracking, while horizontal detail supported finer spatial representation critical for sites like Be'er Sheva with localized dust sources.

iii. **First Layer Height and Emission Schemes:** Lower first-layer heights (10-20 m) improved the models' capacity to capture ground-level concentrations accurately, aligning predictions with near-surface PM10 observations. Emission schemes, especially those incorporating sandblasting (e.g., Marticorena and Bergametti, 1995), supported accurate predictions during intense dust events by modelling larger particle transport effectively.

4. Role of Multi-Model Ensemble (MMM): The multi-model ensemble, which averaged predictions across all models, consistently reduced biases and RMSE values, supporting its effectiveness as a robust forecasting operational approach. Although it does not consistently outperform individual models such as NASA-GEOS, the MMM's reduced error scores make it a reliable option for broader operational applications in the Eastern Mediterranean region (EMR).

5. Approaches and methods of evaluation: The different approaches $O_{PM_{10}}$ and O_{coarse} , different datasets and methods of evaluation have shown that some model's performance can be perceived differently if not properly evaluated by different methods. An example is EMA-REG4, that while being evaluated by performance metrics and BOOT graph methodology it has a very poor performance but on the contingency tables it has one of the best overall scores. Furthermore, swapping from $O_{PM_{10}}$ to O_{coarse} , it shows how sensitive are models to fine materials as they do or do not take them into account changing an underestimation to overestimation and vice versa, while their correlation is mostly decreasing.

This study has limitations, including the absence of recent PM_{2.5} data in some regions and potential changes in model configurations that might affect performance. Despite these limitations, operational models remain reliable for predicting DDS events and provide valuable mitigation insights (Eleftheriou et al., 2023). However, continuous technological advancements are essential to enhance dust prediction capabilities, particularly concerning the duration and intensity of DDS events.

6.2. Future Directions

To advance dust forecasting in the Eastern Mediterranean Region (EMR), targeted model improvements and data integration strategies are essential. Enhancing data assimilation through high-frequency satellite data, such as MODIS AOD, and incorporating localized PM monitoring can improve real-time prediction accuracy, particularly when dynamic field measurements like PM_{2.5} data are included. Importantly, expanding data assimilation to include more in-situ measurements by strengthening monitoring networks, especially in underrepresented source areas, would enhance the accuracy of model predictions and ensure better representation of dust transport from these critical regions. Also, sometimes the models fail producing negative or even zero concentration thus influencing the evaluation technique used to assess their performance. Such results should be redefined to a Limit of Detection (LOD) consistent with observation techniques to provide a more realistic result and the result to be valid in order to be evaluated.

The variability in prediction accuracy across models also highlights the need for improved, site-specific emission parameterizations. Adaptive emission schemes that account for local soil characteristics, vegetation cover, and land use would be especially beneficial for urban-adjacent and desert-proximal sites like Be'er Sheva. Furthermore, results consistently favouring models with finer spatial (0.1° – 0.3°) and vertical resolutions (over 50 layers) suggest that these refinements are crucial for accurately capturing dust dynamics and transboundary events impacting EMR air quality.

Customizing multi-model ensemble strategies by including only the top-performing models for each site could enhance forecast precision while preserving the bias reduction advantages seen in multi-model ensembles. Such an approach leverages specific



model strengths under varied conditions, especially during peak dust events. Expanding model validation efforts to include diverse datasets—such as vertical aerosol profiles, aerosol chemistry, and source-specific particle data—would further refine model parameterizations and improve reliability. Collaborating with regional networks to utilize observational data from North Africa, the Middle East, and EMR could substantially strengthen model robustness.

Future efforts should focus on refining model parameterizations, integrating localized field data, and validating models against diverse datasets to ensure predictive robustness and reliability. Overall, the results underscore the region-specific nature of dust forecasting and highlight model configurations and features that enhance accuracy at individual sites.

Code and Data Availability.

Model forecast data from the SDS-WAS Barcelona Dust Forecast Center (<https://sds-was.aemet.es/>) are openly accessible and were archived during the analysis period. The ground-based measurement data used in this study are available upon request from the authors. This study was performed using custom scripts in MATLAB (more details can be found at <https://www.mathworks.com>) and R (it is available at <https://www.r-project.org>). While no new modelling code was developed in this study, the scripts used for data processing, correlation analysis, and performance evaluation are available upon request from the corresponding author.

Author contributions.

MNKA, PM, and AE contributed to manuscript preparation. MNKA and PY provided the financial support for this study. MNKA also contributed to the investigation, provided resources, and participated in the supervision of the study. PM had the conceptualization, investigation, formal analysis, and supervision of this study. AE and PM performed the formal analysis of the data; AE also prepared the original draft of the manuscript. NK, IK, EV, and CS were responsible for data curation and contributed to manuscript review and editing. PK and PY contributed to writing – review and editing. All co-authors reviewed and approved the final version of the manuscript.

Competing interests.

The contact author has declared that none of the authors has any competing interests.

Acknowledgments

This research work is part of the CiROCCO Project, funded by the European Union (Grant no. 101086497). Views and opinions expressed are, however, those of the author(s) only and do not necessarily reflect those of the European Union or REA. Neither the European Union nor the granting authority can be held responsible for them. Moreover, the authors are grateful for the financial support from the LIFE+-MEDEA Program under Grant Agreement LIFE16 CCA/CY/000041. The authors would also like to express their deepest gratitude to the data providers, the organizations, and their researchers who develop and operate the operational models, as well as to the SDS-WAS portal for the collection and availability of the data.

References

Achilleos, S., Al-Ozairi, E., Alahmad, B., Garshick, E., Neophytou, A. M., Bouhamra, W., Yassin, M. F., and Koutrakis, P.: Acute effects of air pollution on mortality: A 17-year analysis in Kuwait, *Environment International*, 126, 476–483, <https://doi.org/10.1016/j.envint.2019.01.072>, 2019.



- Achilleos, S., Mouzourides, P., Kalivitis, N., Katra, I., Kloog, I., Kouis, P., Middleton, N., Mihalopoulos, N., Neophytou, M., Panayiotou, A., Papatheodorou, S., Savvides, C., Tymvios, F., Vasiliadou, E., Yiallourous, P., and Koutrakis, P.: Spatio-temporal variability of desert dust storms in Eastern Mediterranean (Crete, Cyprus, Israel) between 2006 and 2017 using a uniform methodology, *Science of The Total Environment*, 714, 136693, <https://doi.org/10.1016/j.scitotenv.2020.136693>, 2020.
- Achilleos, S., Michanikou, A., Kouis, P., Papatheodorou, S. I., Panayiotou, A. G., Kinni, P., Mihalopoulos, N., Kalivitis, N., Kouvarakis, G., Galanakis, E., Michailidi, E., Tymvios, F., Chrysanthou, A., Neophytou, M., Mouzourides, P., Savvides, C., Vasiliadou, E., Papasavvas, I., Christophides, T., Nicolaou, R., Avraamides, P., Kang, C.-M., Middleton, N., Koutrakis, P., and Yiallourous, P. K.: Improved indoor air quality during desert dust storms: The impact of the MEDEA exposure-reduction strategies, *Science of The Total Environment*, 863, 160973, <https://doi.org/10.1016/j.scitotenv.2022.160973>, 2023.
- 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.*, 106, 18075–18084, <https://doi.org/10.1029/2000JD900339>, 2001.
- Basart, S., Pérez, C., Nickovic, S., Cuevas, E., and Baldasano, J. M.: Development and evaluation of the BSC-DREAM8b dust regional model over Northern Africa, the Mediterranean and the Middle East, *Tellus B: Chemical and Physical Meteorology*, 64, 18539, <https://doi.org/10.3402/tellusb.v64i0.18539>, 2012.
- Bodenheimer, S., Lensky, I. M., and Dayan, U.: Characterization of Eastern Mediterranean dust storms by area of origin; North Africa vs. Arabian Peninsula, *Atmospheric Environment*, 198, 158–165, <https://doi.org/10.1016/j.atmosenv.2018.10.034>, 2019.
- Çapraz, Ö. and Deniz, A.: Particulate matter (PM₁₀ and PM_{2.5}) concentrations during a Saharan dust episode in Istanbul, *Air Qual Atmos Health*, 14, 109–116, <https://doi.org/10.1007/s11869-020-00917-4>, 2021.
- Chang, J. C. and Hanna, S. R.: Air quality model performance evaluation, *Meteorol Atmos Phys*, 87, <https://doi.org/10.1007/s00703-003-0070-7>, 2004.
- Chen, J. and Hoek, G.: Long-term exposure to PM and all-cause and cause-specific mortality: A systematic review and meta-analysis, *Environment International*, 143, 105974, <https://doi.org/10.1016/j.envint.2020.105974>, 2020.
- Colarco, P., Da Silva, A., Chin, M., and Diehl, T.: Online simulations of global aerosol distributions in the NASA GEOS-4 model and comparisons to satellite and ground-based aerosol optical depth, *J. Geophys. Res.*, 115, 2009JD012820, <https://doi.org/10.1029/2009JD012820>, 2010.
- Eleftheriou, A., Mouzourides, P., Biskos, G., Yiallourous, P., Kumar, P., and Neophytou, M. K.-A.: The challenge of adopting mitigation and adaptation measures for the impacts of sand and dust storms in Eastern Mediterranean Region: a critical review, *Mitig Adapt Strateg Glob Change*, 28, 33, <https://doi.org/10.1007/s11027-023-10070-9>, 2023.
- Flaounas, E., Kotroni, V., Lagouvardos, K., Klose, M., Flamant, C., and Giannaros, T. M.: Sensitivity of the WRF-Chem (V3.6.1) model to different dust emission parametrisation: assessment in the broader Mediterranean region, *Geosci. Model Dev.*, 10, 2925–2945, <https://doi.org/10.5194/gmd-10-2925-2017>, 2017.
- García-Castrillo, G. and E. Terradellas: Evaluation of the dust Forecasts in the Canary Islands, WMO SDS-WAS, Barcelona, 21 pp. SDS-WAS-2017-002, 2017.
- Ginoux, P., Chin, M., Tegen, I., Prospero, J. M., Holben, B., Dubovik, O., and Lin, S.: Sources and distributions of dust aerosols simulated with the GOCART model, *J. Geophys. Res.*, 106, 20255–20273, <https://doi.org/10.1029/2000JD000053>, 2001.
- Holben, B. N., Eck, T. F., Slutsker, I., Tanré, D., Buis, J. P., Setzer, A., Vermote, E., Reagan, J. A., Kaufman, Y. J., Nakajima, T., Lavenue, F., Jankowiak, I., and Smirnov, A.: AERONET—A Federated Instrument Network and Data Archive for Aerosol Characterization, *Remote Sensing of Environment*, 66, 1–16, [https://doi.org/10.1016/S0034-4257\(98\)00031-5](https://doi.org/10.1016/S0034-4257(98)00031-5), 1998.



- Janjić, Z. I.: The Step-Mountain Eta Coordinate Model: Further Developments of the Convection, Viscous Sublayer, and Turbulence Closure Schemes, *Mon. Wea. Rev.*, 122, 927–945, [https://doi.org/10.1175/1520-0493\(1994\)122<0927:TSMECM>2.0.CO;2](https://doi.org/10.1175/1520-0493(1994)122<0927:TSMECM>2.0.CO;2), 1994.
- 695 Jones, B. A.: Dust storms and human well-being, *Resource and Energy Economics*, 72, 101362, <https://doi.org/10.1016/j.reseneeco.2023.101362>, 2023.
- Kazadzis, S., Bais, A., Balis, D., Kouremeti, N., Zempila, M., Arola, A., Giannakaki, E., Amiridis, V., and Kazantzidis, A.: Spatial and temporal UV irradiance and aerosol variability within the area of an OMI satellite pixel, *Atmos. Chem. Phys.*, 9, 4593–4601, <https://doi.org/10.5194/acp-9-4593-2009>, 2009.
- 700 Klose, M., Jorba, O., Gonçalves Ageitos, M., Escibano, J., Dawson, M. L., Obiso, V., Di Tomaso, E., Basart, S., Montané Pinto, G., Macchia, F., Ginoux, P., Guerschman, J., Prigent, C., Huang, Y., Kok, J. F., Miller, R. L., and Pérez García-Pando, C.: Mineral dust cycle in the Multiscale Online Nonhydrostatic Atmosphere Chemistry model (MONARCH) Version 2.0, *Geosci. Model Dev.*, 14, 6403–6444, <https://doi.org/10.5194/gmd-14-6403-2021>, 2021.
- Knippertz, P. and Stuut, J.-B. W. (Eds.): *Mineral Dust: A Key Player in the Earth System*, Springer Netherlands, Dordrecht, <https://doi.org/10.1007/978-94-017-8978-3>, 2014.
- 705 Krasnov, H., Katra, I., and Friger, M.: Increase in dust storm related PM10 concentrations: A time series analysis of 2001–2015, *Environmental Pollution*, 213, 36–42, <https://doi.org/10.1016/j.envpol.2015.10.021>, 2016.
- Lahiri, S. N. (Ed.): *Resampling methods for dependent data*, Springer, New York Berlin Heidelberg, 374 pp., 2010.
- Lorentzou, C., Kouvarakis, G., Kozyrakis, G. V., Kampanis, N. A., Trahanatzi, I., Fraidakis, O., Tzanakis, N., Kanakidou, M., Agouridakis, P., and Notas, G.: Extreme desert dust storms and COPD morbidity on the island of Crete, *COPD*, Volume 14, 1763–1768, <https://doi.org/10.2147/COPD.S208108>, 2019.
- 710 Lu, C.-H., Da Silva, A., Wang, J., Moorthi, S., Chin, M., Colarco, P., Tang, Y., Bhattacharjee, P. S., Chen, S.-P., Chuang, H.-Y., Juang, H.-M. H., McQueen, J., and Iredell, M.: The implementation of NEMS GFS Aerosol Component (NGAC) Version 1.0 for global dust forecasting at NOAA/NCEP, *Geosci. Model Dev.*, 9, 1905–1919, <https://doi.org/10.5194/gmd-9-1905-2016>, 2016.
- 715 Lumet, E., Jaravel, T., Rochoux, M. C., Vermorel, O., and Lacroix, S.: Assessing the Internal Variability of Large-Eddy Simulations for Microscale Pollutant Dispersion Prediction in an Idealized Urban Environment, *Boundary-Layer Meteorol.*, 190, 9, <https://doi.org/10.1007/s10546-023-00853-7>, 2024.
- Lwin, K. S., Tobias, A., Chua, P. L., Yuan, L., Thawonmas, R., Ith, S., Htay, Z. W., Yu, L. S., Yamasaki, L., Roqué, M., Querol, X., Fussell, J. C., Nadeau, K. C., Stafoggia, M., Saliba, N. A., Sheng Ng, C. F., and Hashizume, M.: Effects of Desert Dust and Sandstorms on Human Health: A Scoping Review, *GeoHealth*, 7, e2022GH000728, <https://doi.org/10.1029/2022GH000728>, 2023.
- 720 Manders, A. M. M., Builtjes, P. J. H., Curier, L., Denier Van Der Gon, H. A. C., Hendriks, C., Jonkers, S., Kranenburg, R., Kuenen, J. J. P., Segers, A. J., Timmermans, R. M. A., Visschedijk, A. J. H., Wichink Kruit, R. J., Van Pul, W. A. J., Sauter, F. J., Van Der Swaluw, E., Swart, D. P. J., Douros, J., Eskes, H., Van Meijgaard, E., Van Ulft, B., Van Velthoven, P., Banzhaf, S., Mues, A. C., Stern, R., Fu, G., Lu, S., Heemink, A., Van Velzen, N., and Schaap, M.: Curriculum vitae of the LOTOS–EUROS (v2.0) chemistry transport model, *Geosci. Model Dev.*, 10, 4145–4173, <https://doi.org/10.5194/gmd-10-4145-2017>, 2017.
- 725 Marticorena, B. and Bergametti, G.: Modeling the atmospheric dust cycle: 1. Design of a soil-derived dust emission scheme, *J. Geophys. Res.*, 100, 16415–16430, <https://doi.org/10.1029/95JD00690>, 1995.
- 730



- Mircea, M., D'Isidoro, M., Maurizi, A., Vitali, L., Monforti, F., Zanini, G., and Tampieri, F.: A comprehensive performance evaluation of the air quality model BOLCHEM to reproduce the ozone concentrations over Italy, *Atmospheric Environment*, 42, 1169–1185, <https://doi.org/10.1016/j.atmosenv.2007.10.043>, 2008.
- 735 Morcrette, J. -J., Beljaars, A., Benedetti, A., Jones, L., and Boucher, O.: Sea-salt and dust aerosols in the ECMWF IFS model, *Geophysical Research Letters*, 35, 2008GL036041, <https://doi.org/10.1029/2008GL036041>, 2008.
- Morcrette, J. -J., Boucher, O., Jones, L., Salmond, D., Bechtold, P., Beljaars, A., Benedetti, A., Bonet, A., Kaiser, J. W., Razingzer, M., Schulz, M., Serrar, S., Simmons, A. J., Sofiev, M., Suttie, M., Tompkins, A. M., and Untch, A.: Aerosol analysis and forecast in the European Centre for Medium-Range Weather Forecasts Integrated Forecast System: Forward modeling, *J. Geophys. Res.*, 114, 2008JD011235, <https://doi.org/10.1029/2008JD011235>, 2009.
- 740 Mouzourides, P., Kumar, P., and Neophytou, M. K.-A.: Assessment of long-term measurements of particulate matter and gaseous pollutants in South-East Mediterranean, *Atmospheric Environment*, 107, 148–165, <https://doi.org/10.1016/j.atmosenv.2015.02.031>, 2015.
- NASA Publications: MODIS, <https://terra.nasa.gov/about/terra-instruments/modis> (last access 27. February, 2025), 2025.
- Nickovic, S., Kallos, G., Papadopoulos, A., and Kakaliagou, O.: A model for prediction of desert dust cycle in the atmosphere, *J. Geophys. Res.*, 106, 18113–18129, <https://doi.org/10.1029/2000JD900794>, 2001.
- 745 Nickovic, S., Cvetkovic, B., Madonna, F., Rosoldi, M., Pejanovic, G., Petkovic, S., and Nikolic, J.: Cloud ice caused by atmospheric mineral dust – Part 1: Parameterization of ice nuclei concentration in the NMME-DREAM model, *Atmos. Chem. Phys.*, 16, 11367–11378, <https://doi.org/10.5194/acp-16-11367-2016>, 2016.
- Pérez, C., Nickovic, S., Pejanovic, G., Baldasano, J. M., and Özsoy, E.: Interactive dust-radiation modeling: A step to improve weather forecasts, *J. Geophys. Res.*, 111, 2005JD006717, <https://doi.org/10.1029/2005JD006717>, 2006.
- 750 Pérez, C., Haustein, K., Janjic, Z., Jorba, O., Huneeus, N., Baldasano, J. M., Black, T., Basart, S., Nickovic, S., Miller, R. L., Perlwitz, J. P., Schulz, M., and Thomson, M.: Atmospheric dust modeling from meso to global scales with the online NMMB/BSC-Dust model – Part 1: Model description, annual simulations and evaluation, <https://doi.org/10.5194/acpd-11-17551-2011>, 22 June 2011.
- 755 Querol, X., Pey, J., Pandolfi, M., Alastuey, A., Cusack, M., Pérez, N., Moreno, T., Viana, M., Mihalopoulos, N., Kallos, G., and Kleanthous, S.: African dust contributions to mean ambient PM₁₀ mass-levels across the Mediterranean Basin, *Atmospheric Environment*, 43, 4266–4277, <https://doi.org/10.1016/j.atmosenv.2009.06.013>, 2009.
- Rémy, S., Kipling, Z., Huijnen, V., Flemming, J., Nabat, P., Michou, M., Ades, M., Engelen, R., and Peuch, V.-H.: Description and evaluation of the tropospheric aerosol scheme in the Integrated Forecasting System (IFS-AER, cycle 47R1) of ECMWF, *Geosci. Model Dev.*, 15, 4881–4912, <https://doi.org/10.5194/gmd-15-4881-2022>, 2022.
- 760 Shao, Y., Raupach, M. R., and Findlater, P. A.: Effect of saltation bombardment on the entrainment of dust by wind, *J. Geophys. Res.*, 98, 12719–12726, <https://doi.org/10.1029/93JD00396>, 1993.
- Sofiev, M., Vira, J., Kouznetsov, R., Prank, M., Soares, J., and Genikhovich, E.: Construction of the SILAM Eulerian atmospheric dispersion model based on the advection algorithm of Michael Galperin, *Geosci. Model Dev.*, 8, 3497–3522, <https://doi.org/10.5194/gmd-8-3497-2015>, 2015.
- 765 Solomos, S., Spyrou, C., Barreto, A., Rodríguez, S., González, Y., Neophytou, M. K. A., Mouzourides, P., Bartsotas, N. S., Kalogeri, C., Nickovic, S., Vukovic Vimic, A., Vujadinovic Mandic, M., Pejanovic, G., Cvetkovic, B., Amiridis, V., Sykioti, O., Gkikas, A., and Zerefos, C.: The Development of METAL-WRF Regional Model for the Description of Dust Mineralogy in the Atmosphere, *Atmosphere*, 14, 1615, <https://doi.org/10.3390/atmos14111615>, 2023.



- 770 United Nations Publications:
<https://press.un.org/en/2024/ga12613.doc.htm#:~:text=The%20193%2Dmember%20organ%20first,fight%20against%20thos,e%20meteorological%20phenomena> (last access 27. February, 2025), 2024.
- Tegen, I. and Fung, I.: Modeling of mineral dust in the atmosphere: Sources, transport, and optical thickness, *J. Geophys. Res.*, 99, 22897–22914, <https://doi.org/10.1029/94JD01928>, 1994.
- 775 Triantafyllou, E., Diapouli, E., Korras-Carraca, M. B., Manousakas, M., Psanis, C., Floutsi, A. A., Spyrou, C., Eleftheriadis, K., and Biskos, G.: Contribution of locally-produced and transported air pollution to particulate matter in a small insular coastal city, *Atmospheric Pollution Research*, 11, 667–678, <https://doi.org/10.1016/j.apr.2019.12.015>, 2020.
- Tsiflikiotou, M. A., Kostenidou, E., Papanastasiou, D. K., Patoulas, D., Zampas, P., Paraskevopoulou, D., Diapouli, E., Kaltsonoudis, C., Florou, K., Bougiatioti, A., Stavroulas, I., Theodosi, C., Kouvarakis, G., Vasilatou, V., Siakavaras, D.,
780 Biskos, G., Pilinis, C., Eleftheriadis, K., Gerasopoulos, E., Mihalopoulos, N., and Pandis, S. N.: Summertime particulate matter and its composition in Greece, *Atmospheric Environment*, 213, 597–607, <https://doi.org/10.1016/j.atmosenv.2019.06.013>, 2019.
- Vratolis, S., Gini, M. I., Bezantakos, S., Stavroulas, I., Kalivitis, N., Kostenidou, E., Louvaris, E., Siakavaras, D., Biskos, G., Mihalopoulos, N., Pandis, S. N., Pilinis, C., Papayannis, A., and Eleftheriadis, K.: Particle number size distribution statistics
785 at City-Centre Urban Background, urban background, and remote stations in Greece during summer, *Atmospheric Environment*, 213, 711–726, <https://doi.org/10.1016/j.atmosenv.2019.05.064>, 2019.
- Yang, L., Fang, S., Zhuang, S., Chen, Y., Li, X., and Zhang, Q.: Atmospheric ¹³⁷Cs dispersion following the Fukushima Daiichi nuclear accident: Local-scale simulations using CALMET and LAPMOD, *Annals of Nuclear Energy*, 195, 110137, <https://doi.org/10.1016/j.anucene.2023.110137>, 2024.
- 790 Zakey, A. S., Solmon, F., and Giorgi, F.: Implementation and testing of a desert dust module in a regional climate model, *Atmos. Chem. Phys.*, 6, 4687–4704, <https://doi.org/10.5194/acp-6-4687-2006>, 2006.
- Zender, C. S., Bian, H., and Newman, D.: Mineral Dust Entrainment and Deposition (DEAD) model: Description and 1990s dust climatology, *J. Geophys. Res.*, 108, 2002JD002775, <https://doi.org/10.1029/2002JD002775>, 2003.
- Zhang, L., Montuoro, R., McKeen, S. A., Baker, B., Bhattacharjee, P. S., Grell, G. A., Henderson, J., Pan, L., Frost, G. J.,
795 McQueen, J., Saylor, R., Li, H., Ahmadov, R., Wang, J., Stajner, I., Kondragunta, S., Zhang, X., and Li, F.: Development and evaluation of the Aerosol Forecast Member in the National Center for Environment Prediction (NCEP)’s Global Ensemble Forecast System (GEFS-Aerosols v1), *Geosci. Model Dev.*, 15, 5337–5369, <https://doi.org/10.5194/gmd-15-5337-2022>, 2022.