

**IDEQ Air Quality
Machine Learning Forecast System
Version 2**

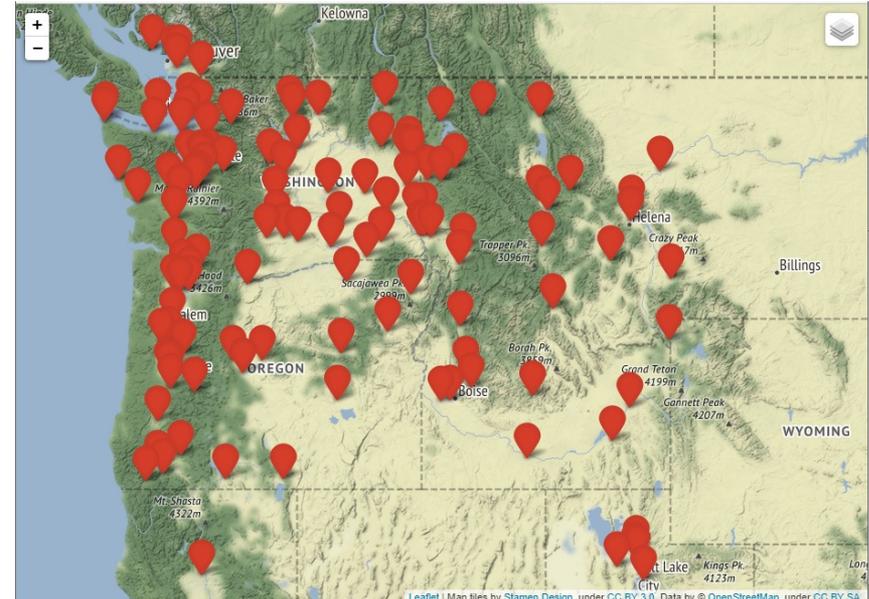
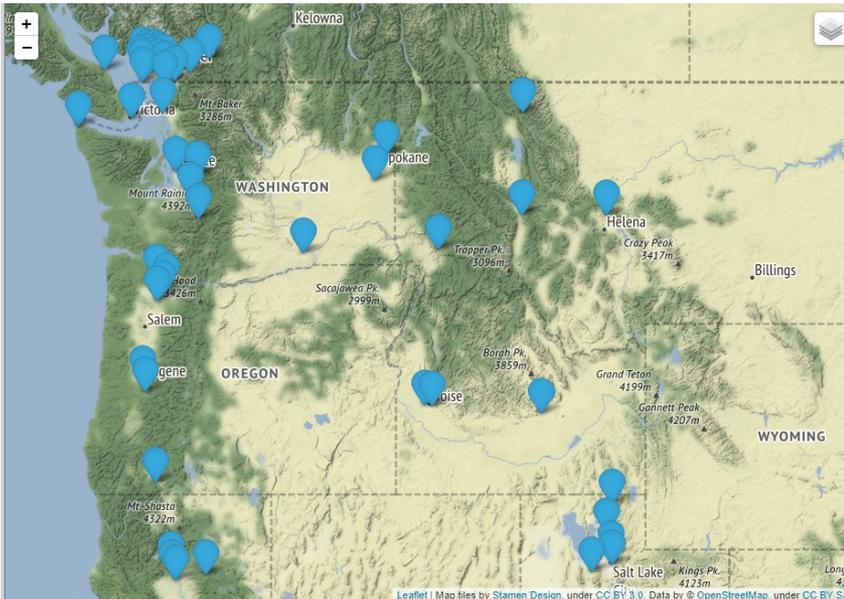
NW-AIRQUEST 2022 Annual Meeting
06/29/2022

Overview of Forecast System

Spatial Coverage

O3 Sites (50, 4 in Idaho)

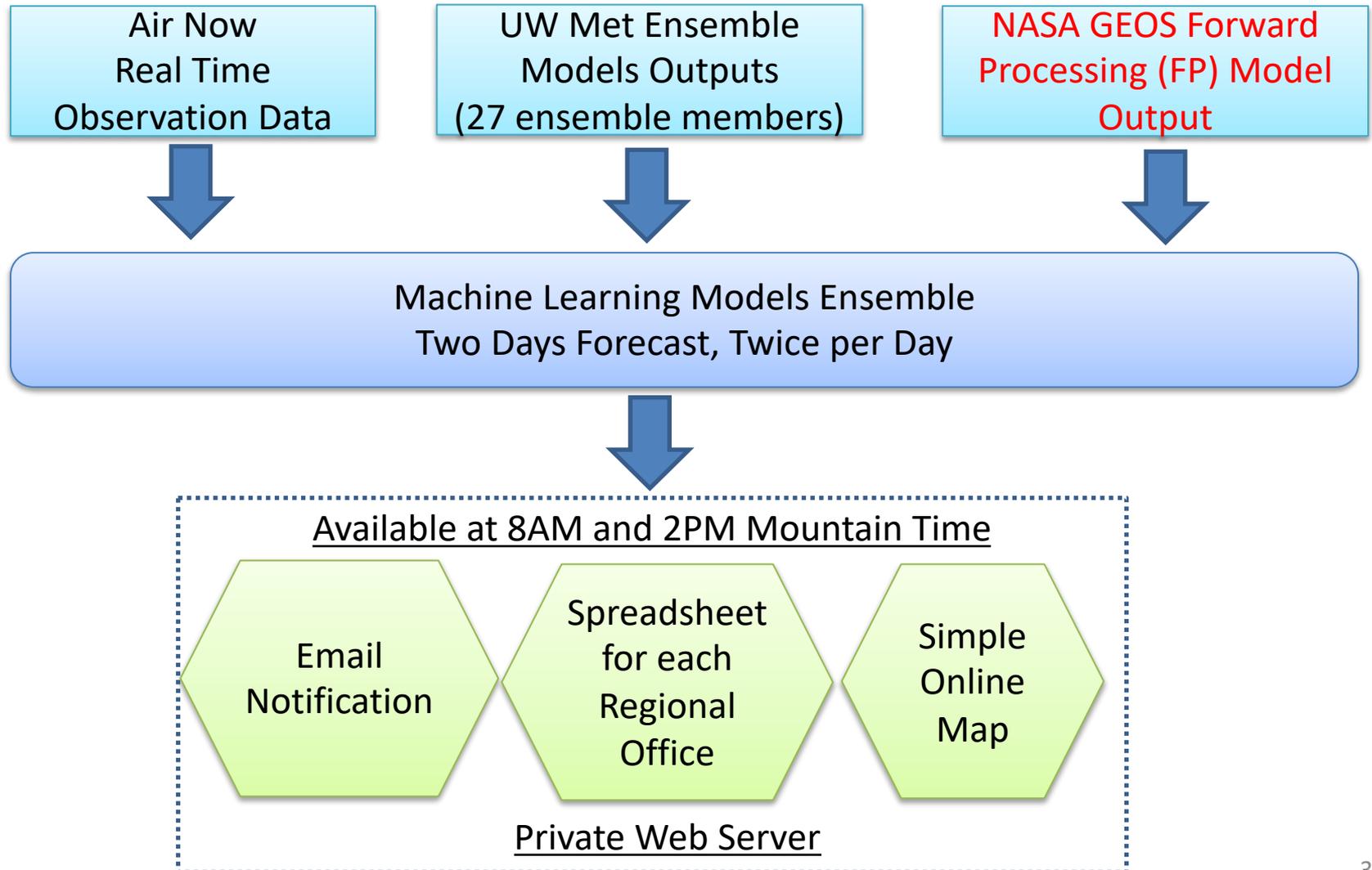
PM2.5 Sites (130, 21 in Idaho)



Site Selection Criteria: Having at least 3 years of data and continue to operate in the future

Overview of Forecast System

Data Process – In and Out



UW Meteorology Models

- Currently 27 ensemble members
- 4 km resolution
- Initialized at 00Z and 12Z
- Each forecasts 72 hours

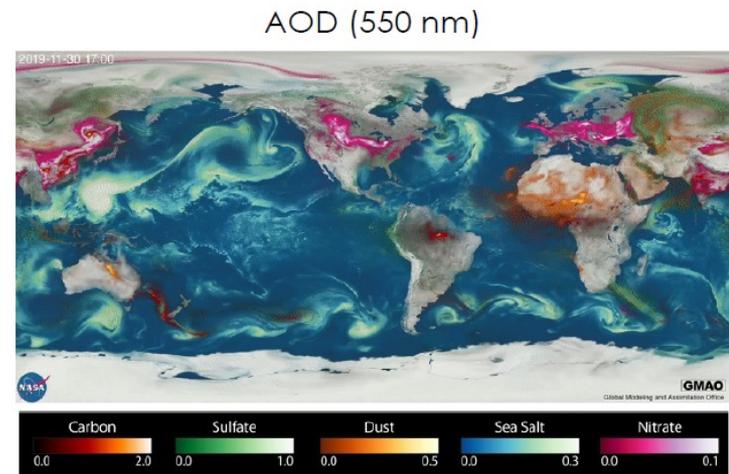
- Time of arriving
 - 00Z forecast : ~09 - 12 MST
 - 12Z forecast : ~21 - 00 MST
- Utilization
 - Morning forecast : using previous day 12Z Met forecast
 - Afternoon forecast : using current day 00Z Met forecast

NASA GEOS-FP

GEOS FP

https://gmao.gsfc.nasa.gov/weather_prediction/

- GEOS FP analyses and forecasts support NASA field campaigns and provide a testbed for assimilation and forecast development
- Publicly available
- Includes weather, aerosols, and carbon monoxide (CO) on the same spatial scale
- State of the science forecast system – model physics or observing system updated every 6-12 months
 - Not suitable for trend analyses
- Meteorology used to drive chemistry models:
 - GEOSChem, Whole Atmosphere Community Climate Model (WACCM)
- When using FP meteorology fields to drive another model, must ensure your simulation does not span an update
 - [GMAO NRT Product Page](#) has updated details and dates



<https://svs.gsfc.nasa.gov/31100>

NASA GEOS-FP Output

GEOS Output Quick Guide

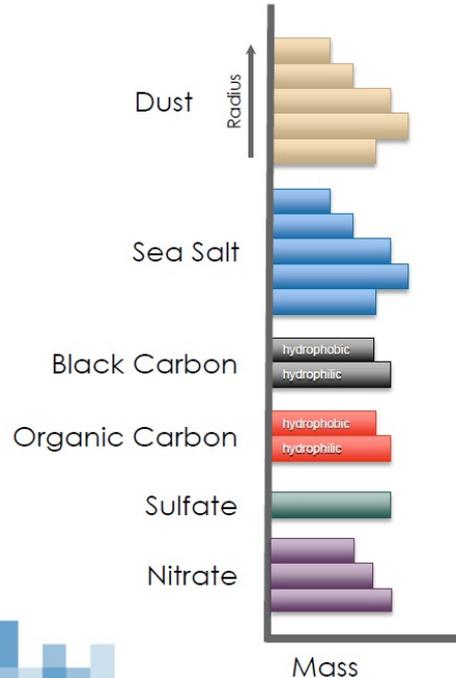
	GEOS FP
Type	Analysis + Forecast
Domain	Global
Spatial Resolution	Simulation: ~12 km Output: ~25 km (0.25°x0.312°)
Temporal Resolution	2-D data: Hourly 3-D data: Every 3 h
Vertical Levels	72 (near surface-0.1 hPa)
Output Available	Analysis: 2014 – Present Forecast: ~21 days
Initialization	Daily 10-day forecast at 00Z Daily 5-day forecast at 12Z
Data Assimilation	Yes
File Specification Doc	https://gmao.gsfc.nasa.gov/pubs/docs/Lucchesi1203.pdf *



GOCART

GOCART in GEOS

- Goddard Chemistry, Aerosol, Radiation and Transport Model (GOCART, Chin et al. 2002, Colarco et al. 2010)
- Sources and sinks for 6 non-interactive species
- Radiatively active



Wind and topographic sources, 5 mass bins

Wind-driven source, 5 mass bins

Anthropogenic and wildfire sources, mass hydrophobic & hydrophilic

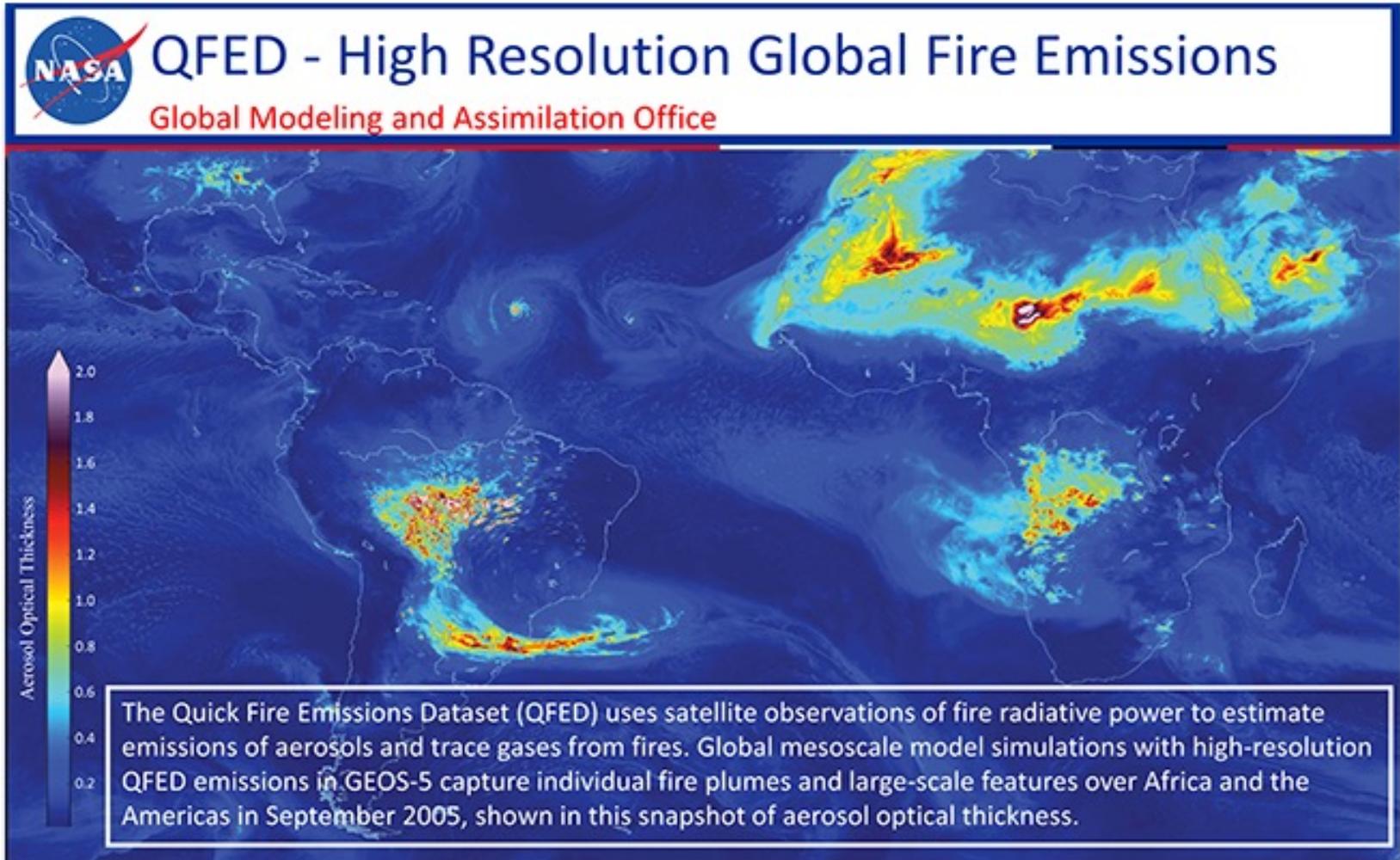
Anthropogenic, biogenic, and fire sources, mass hydrophobic and hydrophilic

Anthropogenic, wildfire, and volcanic

Anthropogenic and wildfire sources



QFED - High Resolution Global Fire Emissions



https://gmao.gsfc.nasa.gov/research/science_snapshots/global_fire_emissions.php

O3 and PM2.5 Forecast Inputs Used

17 Inputs :

- Previous 24th hour O3/PM2.5
- T, P, RH, PBLH, WS, WDIR
- O3/PM2.5 hourly_mean (grouped by Month, Weekday/end, Hour)
- GOES FP PM2.5 Species: Dust, Organic Carbon, etc.
- Month (1-12), Weekday (0-1), Hour (0-23)

O3 and PM2.5 Forecast

Machine Learning Models Used

Neural Network Models

- Dense Neural Network Model
- 1D Convolutional Neural Network Model
- Recurrent Neural Network (LSTM)

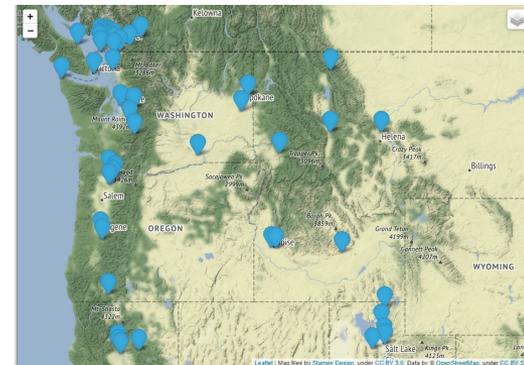
Decision Tree Based Models

- Gradient Boosting
- Random Forest
- Boosted Random Forest

Baseline Model : Persistence

Wildfire Machine Learning Models

- Combine all sites together
 - Get adequate data set to train models
- Introduce site characteristic inputs
 - Elevation, Solar irradiance, Population, Land Use, Terrain Feature
- Limit to wildfire impact period and daytime
- Build models for different concentration ranges
 - Lessen Data imbalance issue
- Combine base model and wildfire models in an escalating fashion

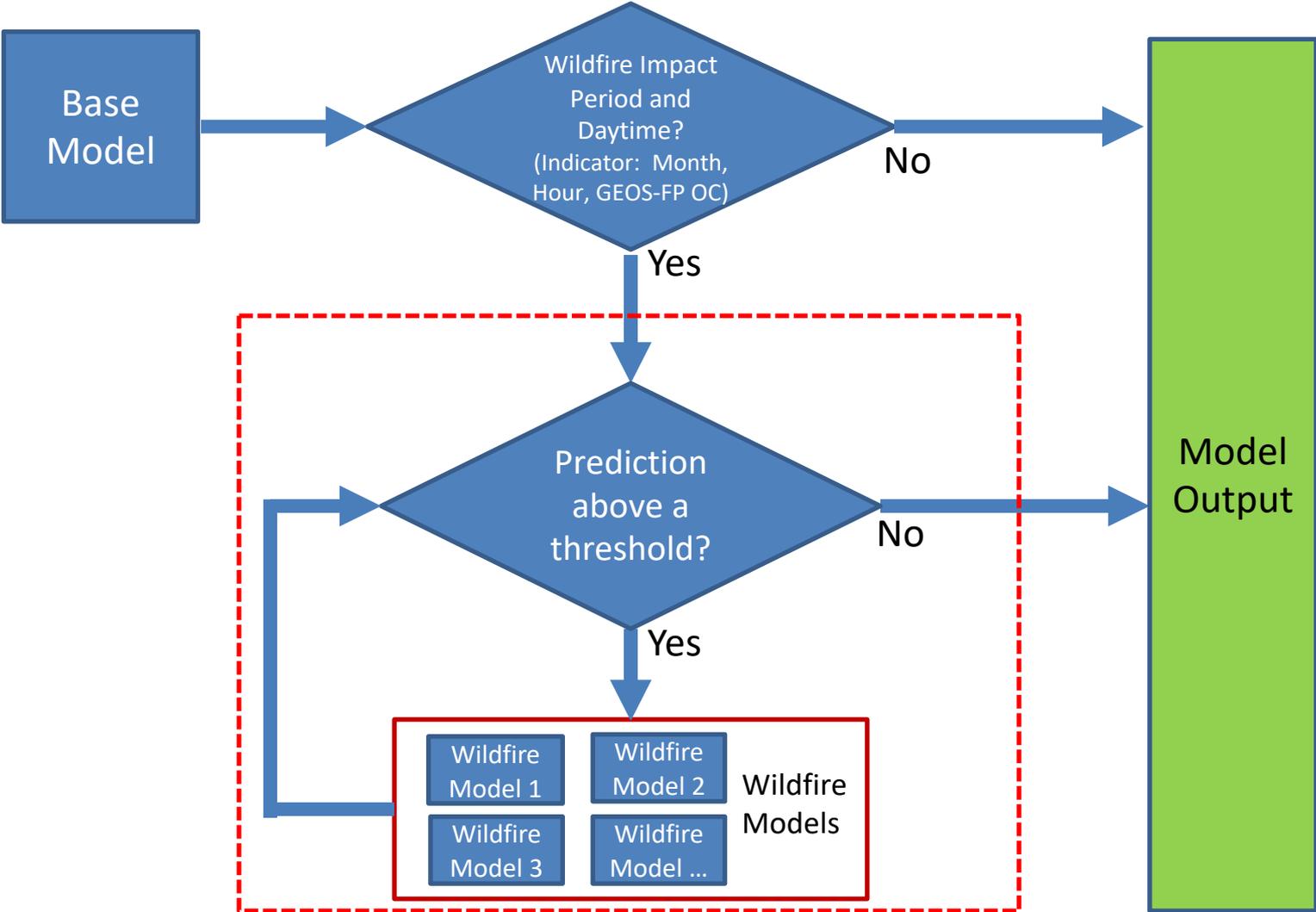


Wildfire Model Stairs



- Wildfire impact period
- Day time
- Upper end of concentration
- Lessen Data imbalance issue

Wildfire Enhanced Model Flow



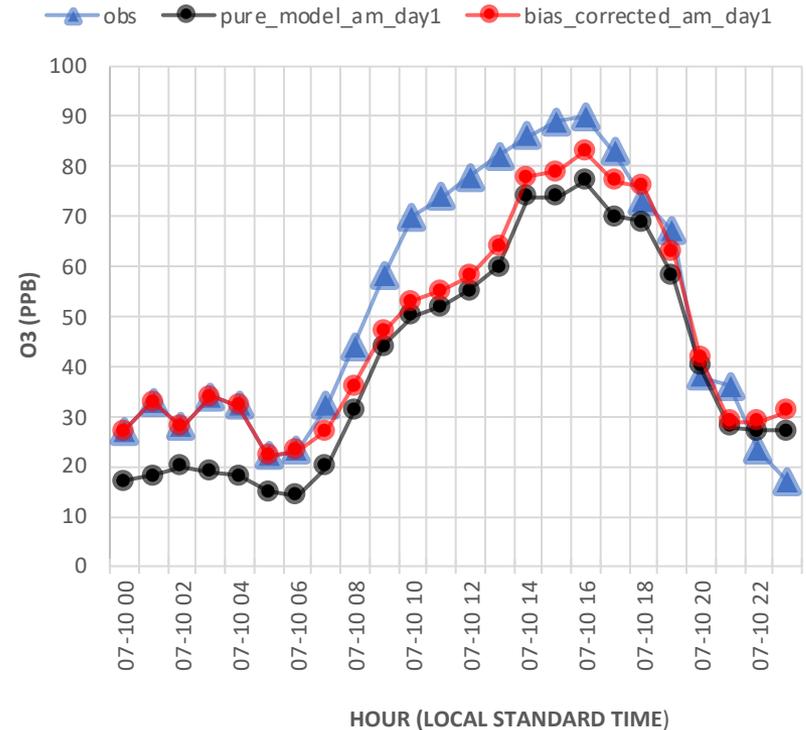
Ensemble

- Three layers of Ensemble
- Ensemble already employed in some Machine Learning Models, such as random forest
- Ensemble of Machine Learning Models:
 - Produce the final model output from multiple machine learning models for one set of input data
- Ensemble of Meteorology Models
 - The distribution of prediction

Bias Correction for Day 1

- Replace prediction with known observation
- Apply bias correction to directly following several hours based on known bias of previous 3 hours
- Apply bias correction for the rest hours based on previous day's 3-hour window hourly bias

O3 Hourly Concentration
Meridian (AQSID 160010010)
Forecast at 07/10/2021 Morning



Annual Maintenance

Retraining Model Annually

- Why
 - Machine Learning Model needs a lot of data
 - Model Performance may drift overtime
- When and How
 - Each January
 - Using previous year's data as test data set and free up part of original test data into training data set
 - Tune up model parameters

Model Performance
at
St. Lukes Meridian
for
03 Models Used in Year 2021

St. Lukes Meridian O3 Site

Year 2021 Model Performance (1)

Daily Regression Performance Metrics

Forecast	max_error	mean_absolute_error	mean_squared_error	normalized_mean_bias	normalized_mean_error	r2_score	root_mean_squared_error
persistence	37	5.69	62.29	0	0.13	0.66	7.89
pure_model_am_day2	22	5.52	47.41	-0.06	0.13	0.74	6.89
pure_model_pm_day2	21	5.27	42.93	-0.05	0.12	0.76	6.55
pure_model_am_day1	20	5.15	41.14	-0.05	0.12	0.77	6.41
pure_model_pm_day1	20	5.11	40.21	-0.06	0.12	0.78	6.34
bias_corrected_am_day2	22	5.23	43.44	-0.05	0.12	0.76	6.59
bias_corrected_pm_day2	23	4.91	40.36	-0.04	0.12	0.78	6.35
bias_corrected_am_day1	25	4.18	31.03	-0.03	0.1	0.83	5.57
bias_corrected_pm_day1	13	2.49	11.51	-0.01	0.06	0.94	3.39

St. Lukes Meridian O3 Site

Year 2021 Model Performance (2)

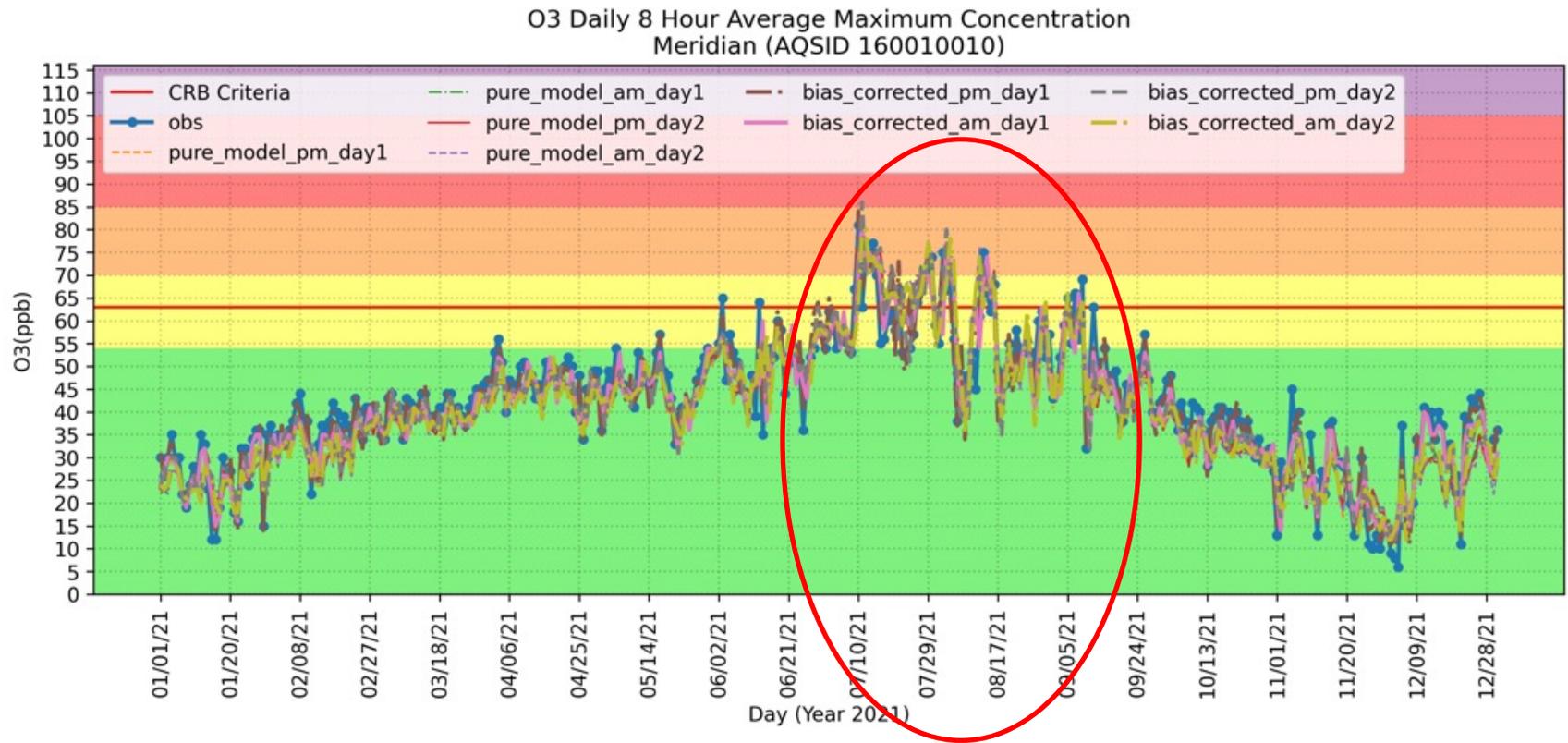
AQI Classification Metrics for All AQI Class

Forecast	AQI_class	accuracy	HSS	KSS
persistence	ALL	0.86	0.5	0.5
pure_model_am_day2	ALL	0.89	0.61	0.62
pure_model_pm_day2	ALL	0.9	0.63	0.63
pure_model_am_day1	ALL	0.9	0.65	0.65
pure_model_pm_day1	ALL	0.9	0.65	0.64
bias_corrected_am_day2	ALL	0.89	0.61	0.61
bias_corrected_pm_day2	ALL	0.9	0.64	0.64
bias_corrected_am_day1	ALL	0.9	0.62	0.62
bias_corrected_pm_day1	ALL	0.92	0.72	0.74

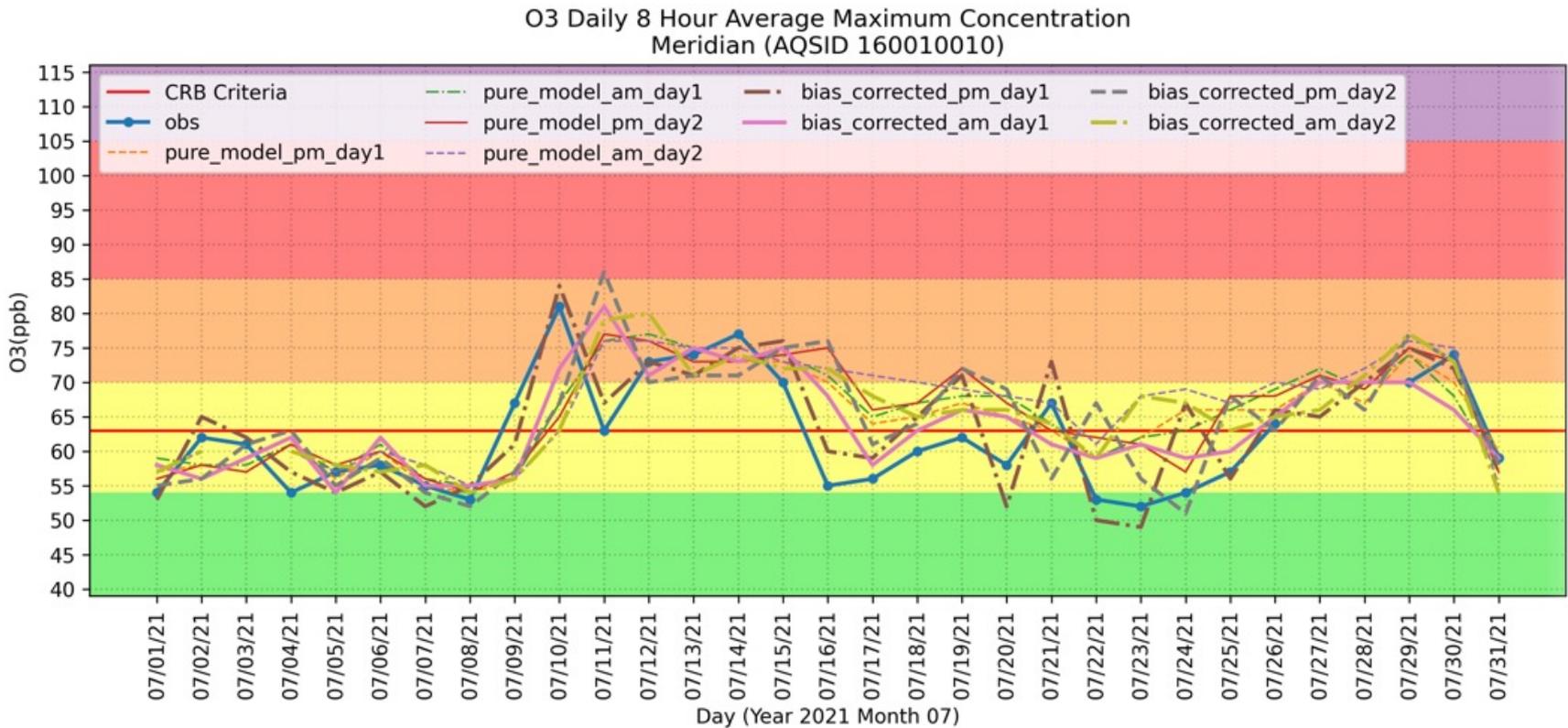
AQI Classification Metrics for AQI Class 2 (Yellow)

Forecast	AQI_class	precision	recall	f1-score	support
persistence	AQI class 2	0.5	0.5	0.5	48
pure_model_am_day2	AQI class 2	0.62	0.55	0.58	47
pure_model_pm_day2	AQI class 2	0.67	0.54	0.6	48
pure_model_am_day1	AQI class 2	0.65	0.62	0.64	48
pure_model_pm_day1	AQI class 2	0.65	0.62	0.64	48
bias_corrected_am_day2	AQI class 2	0.61	0.57	0.59	47
bias_corrected_pm_day2	AQI class 2	0.66	0.6	0.63	48
bias_corrected_am_day1	AQI class 2	0.62	0.6	0.61	48
bias_corrected_pm_day1	AQI class 2	0.7	0.73	0.71	48

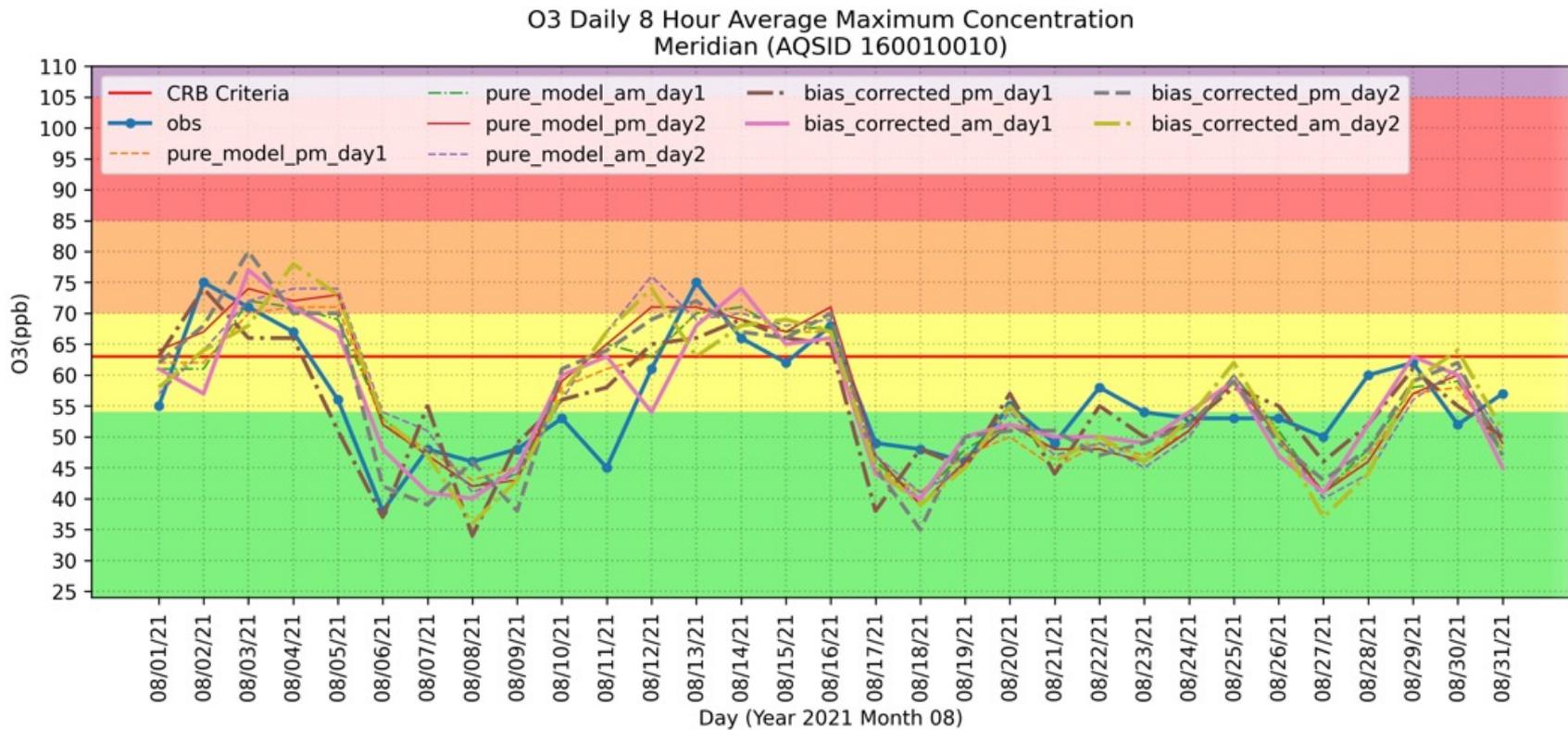
St. Lukes Meridian O3 Site Year 2021 Time Series



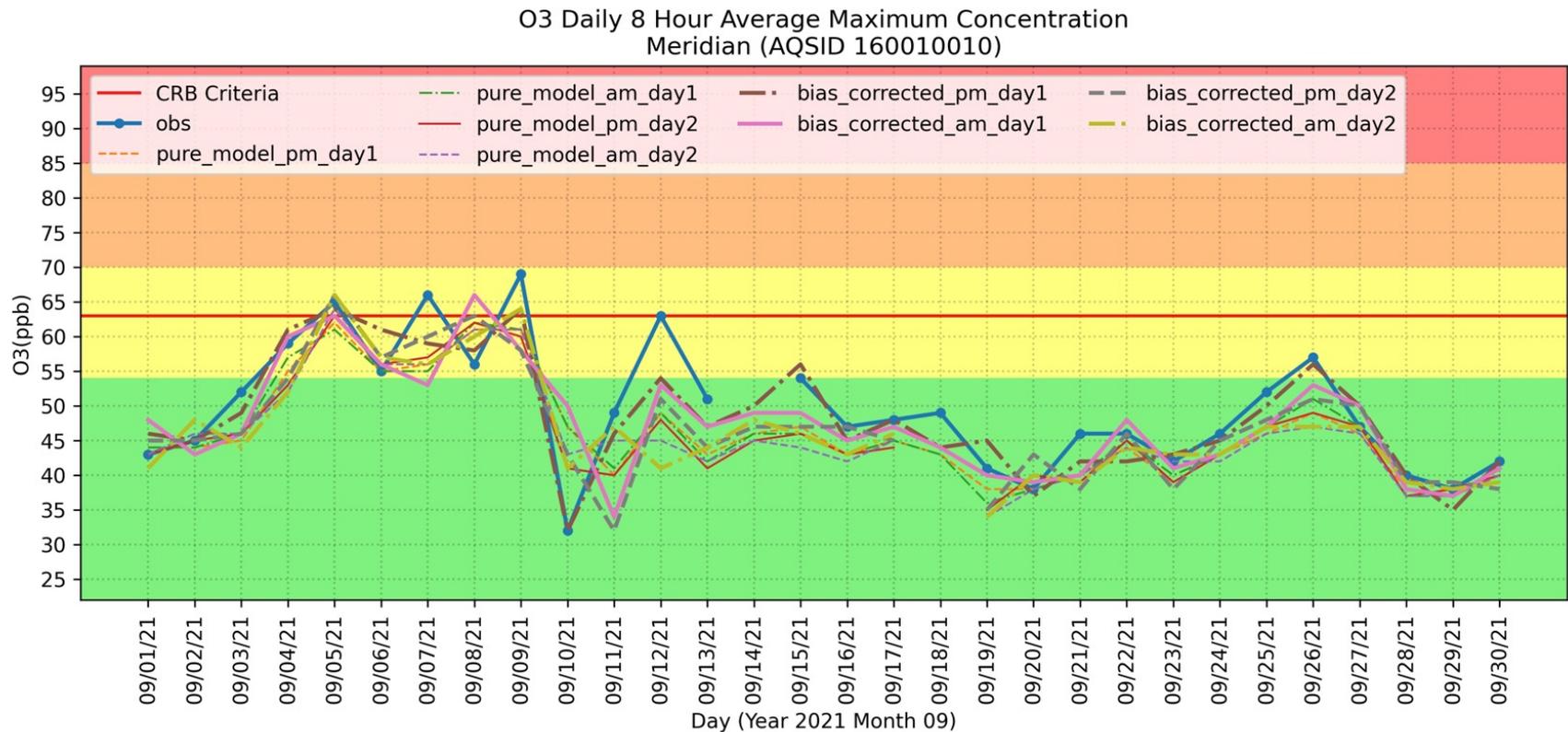
St. Lukes Meridian O3 Site July 2021 Time Series



St. Lukes Meridian O3 Site August 2021 Time Series

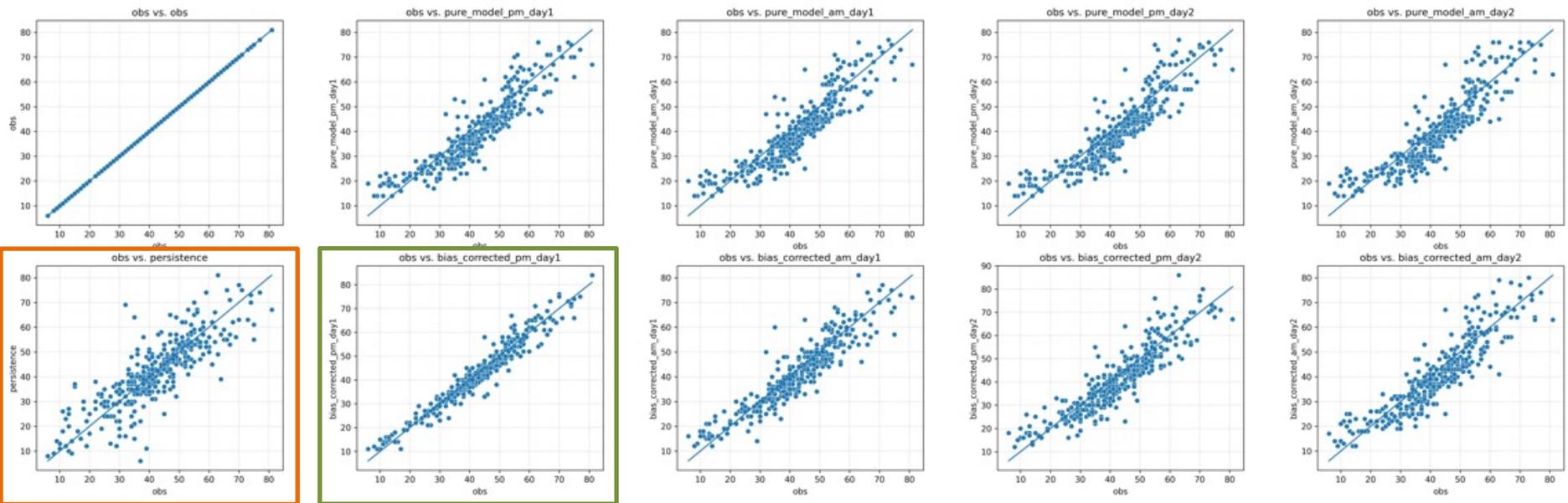


St. Lukes Meridian O3 Site September 2021 Time Series



St. Lukes Meridian O3 Site Year 2021 Daily Scatter Plots

O3 Daily 8 Hour Average Maximum Concentration Scatter Plots
Meridian (AQSID 160010010)
Year 2021



Model Performance
at
St. Lukes Meridian
for
PM2.5 Models

St. Lukes Meridian PM2.5 Site Year 2021 Model Performance (1)

Daily Regression Performance Metrics

Forecast	max_error	mean_absolute_error	mean_squared_error	normalized_mean_bias	normalized_mean_error	r2_score	root_mean_squared_error
persistence	36.1	3.59	35.12	0	0.38	0.62	5.93
pure_model_am_day2	43.9	2.62	20.35	-0.05	0.27	0.78	4.51
pure_model_pm_day2	42.9	2.51	19.33	-0.05	0.26	0.79	4.4
pure_model_am_day1	39.9	2.41	17.59	-0.05	0.25	0.81	4.19
pure_model_pm_day1	39.9	2.41	17.4	-0.05	0.25	0.81	4.17
bias_corrected_am_day2	40.9	2.66	20.19	-0.05	0.28	0.78	4.49
bias_corrected_pm_day2	38.8	2.54	19.26	-0.04	0.27	0.79	4.39
bias_corrected_am_day1	19	1.68	7.94	-0.02	0.18	0.91	2.82
bias_corrected_pm_day1	15.5	1.25	4.56	0	0.13	0.95	2.14

St. Lukes Meridian PM2.5 Site

Year 2021 Model Performance (2)

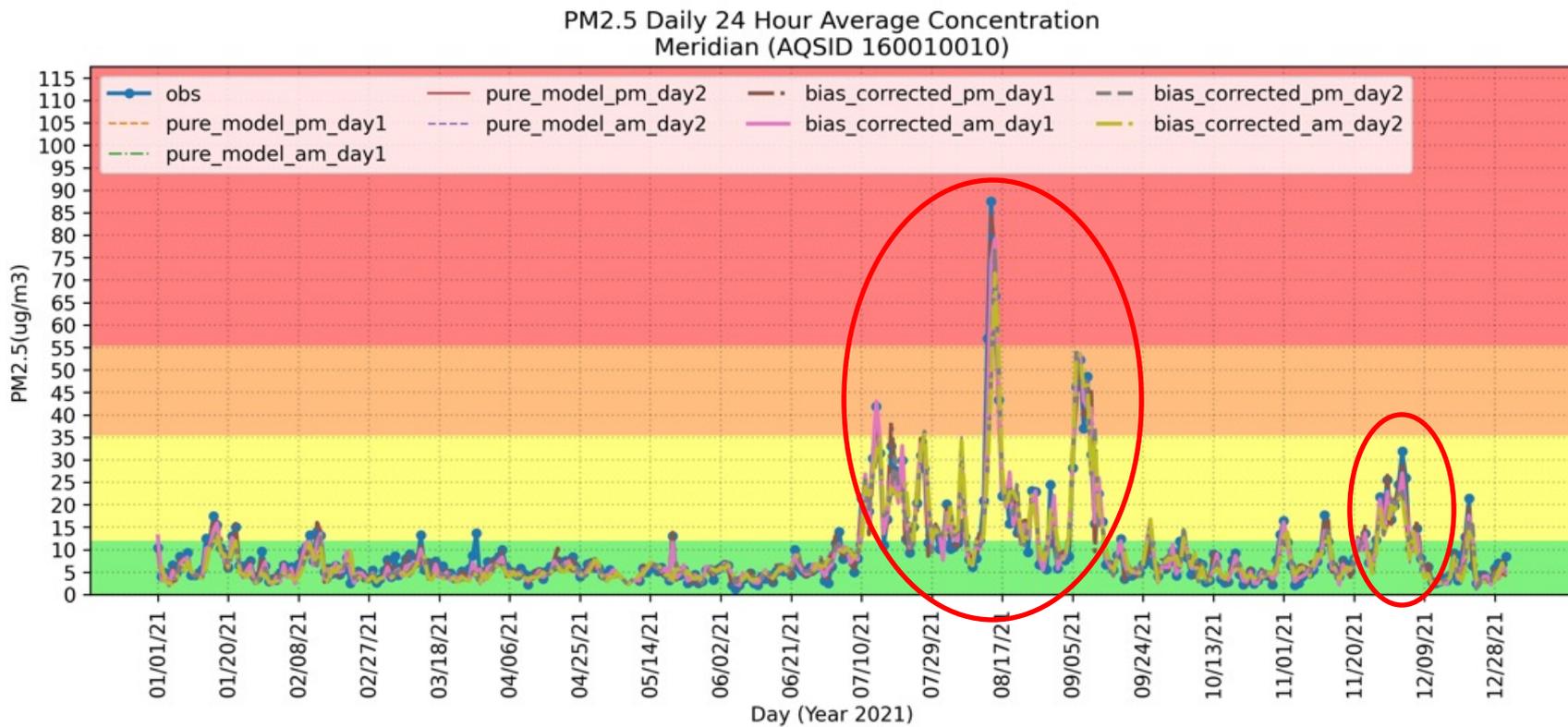
AQI Classification Metrics for All AQI Class

Forecast	AQI_class	accuracy	HSS	KSS
persistence	ALL	0.83	0.5	0.51
pure_model_am_day2	ALL	0.89	0.67	0.65
pure_model_pm_day2	ALL	0.89	0.67	0.66
pure_model_am_day1	ALL	0.89	0.67	0.66
pure_model_pm_day1	ALL	0.89	0.66	0.64
bias_corrected_am_day2	ALL	0.89	0.66	0.65
bias_corrected_pm_day2	ALL	0.88	0.65	0.64
bias_corrected_am_day1	ALL	0.91	0.74	0.74
bias_corrected_pm_day1	ALL	0.93	0.8	0.8

AQI Classification Metrics for AQI Class 2 (Yellow)

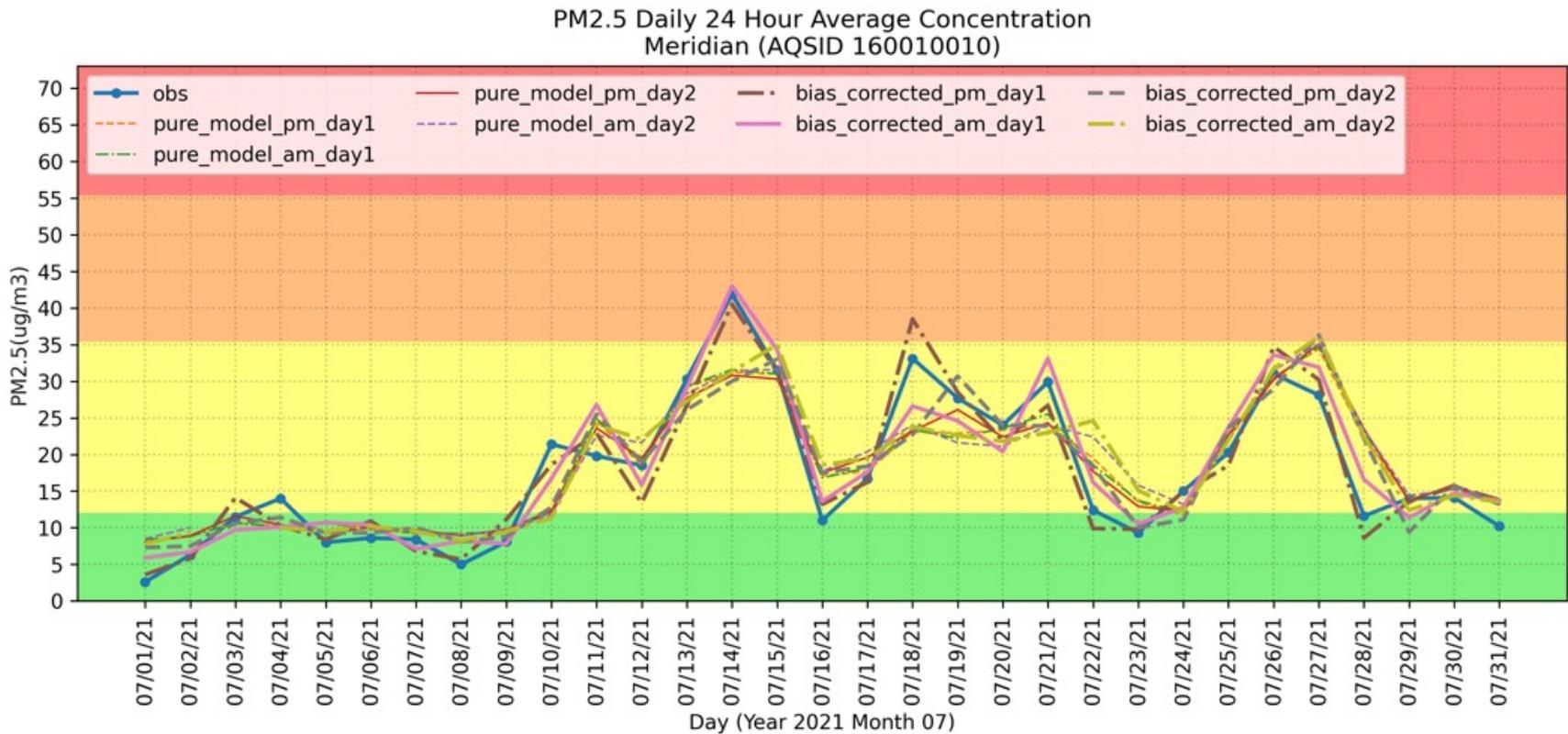
Forecast	AQI_class	precision	recall	f1-score	support
persistence	AQI class 2	0.54	0.56	0.55	66
pure_model_am_day2	AQI class 2	0.71	0.69	0.7	65
pure_model_pm_day2	AQI class 2	0.7	0.71	0.7	65
pure_model_am_day1	AQI class 2	0.7	0.71	0.71	66
pure_model_pm_day1	AQI class 2	0.7	0.68	0.69	66
bias_corrected_am_day2	AQI class 2	0.71	0.68	0.69	65
bias_corrected_pm_day2	AQI class 2	0.69	0.66	0.68	65
bias_corrected_am_day1	AQI class 2	0.77	0.76	0.76	66
bias_corrected_pm_day1	AQI class 2	0.83	0.8	0.82	66

St. Lukes Meridian PM2.5 Site Year 2021 Time Series



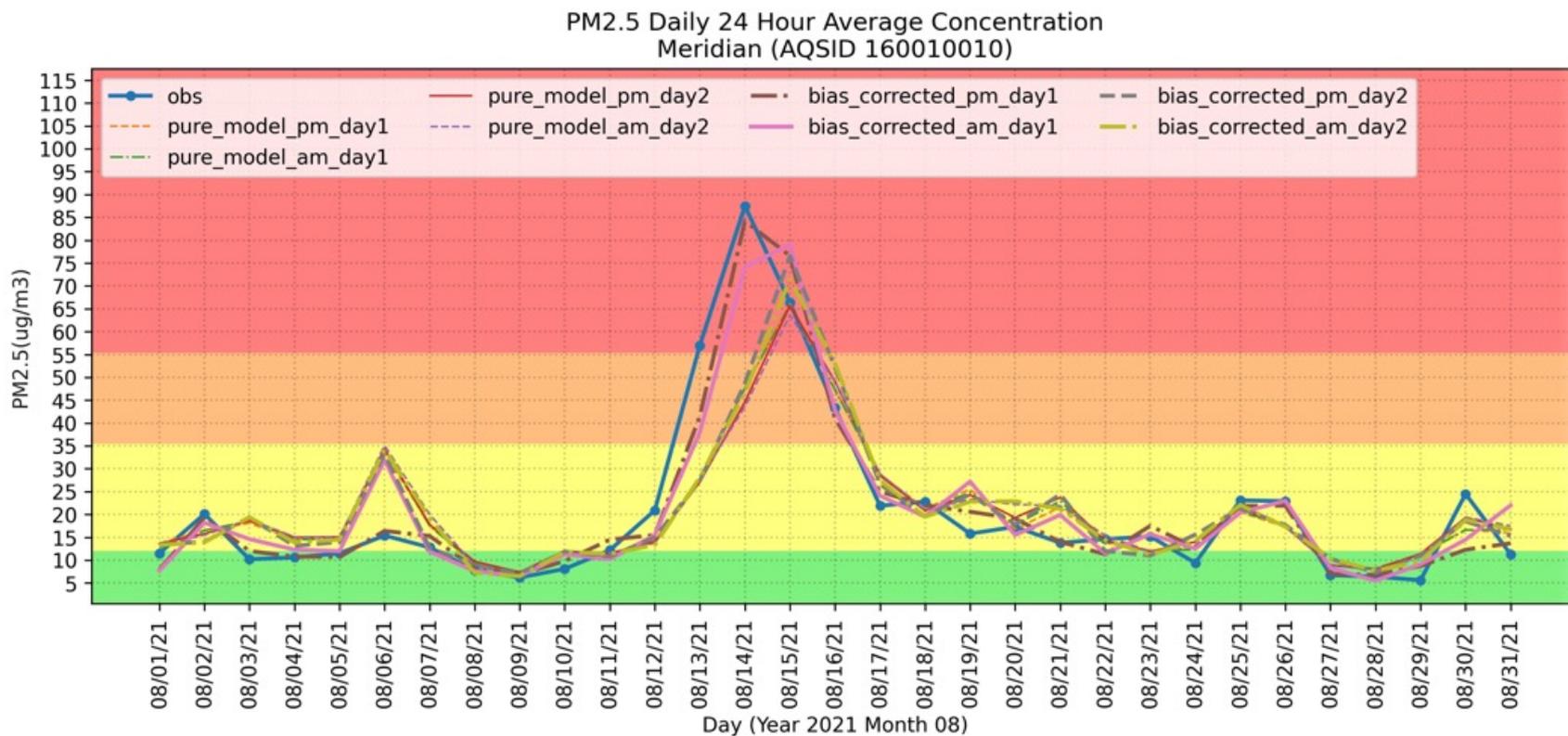
St. Lukes Meridian PM2.5 Site

July 2021 Time Series

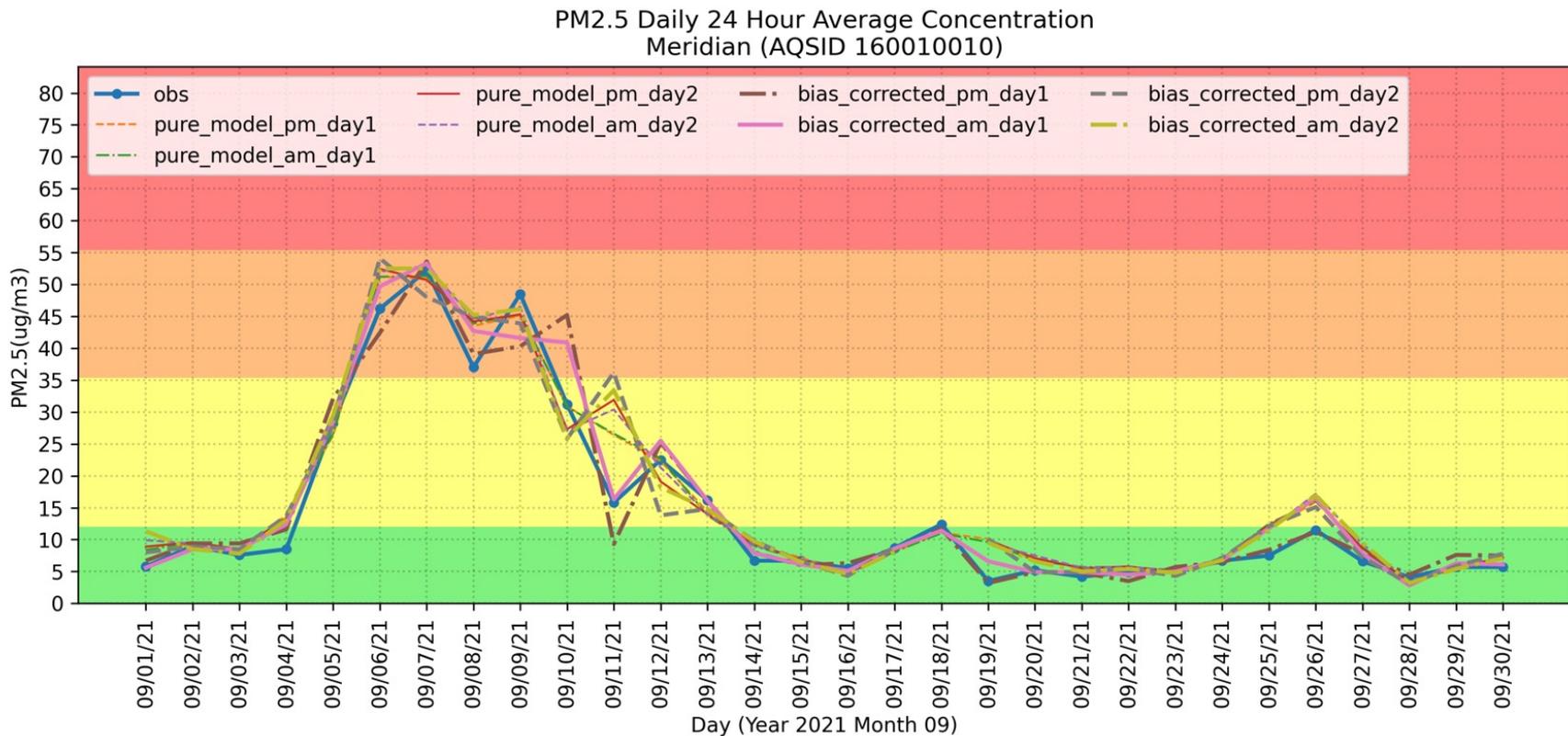


St. Lukes Meridian PM2.5 Site

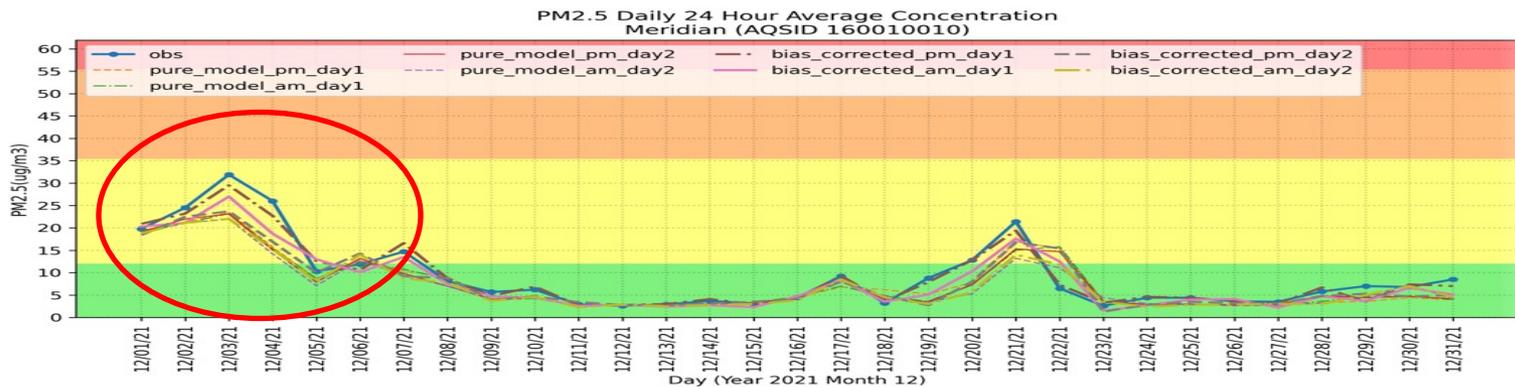
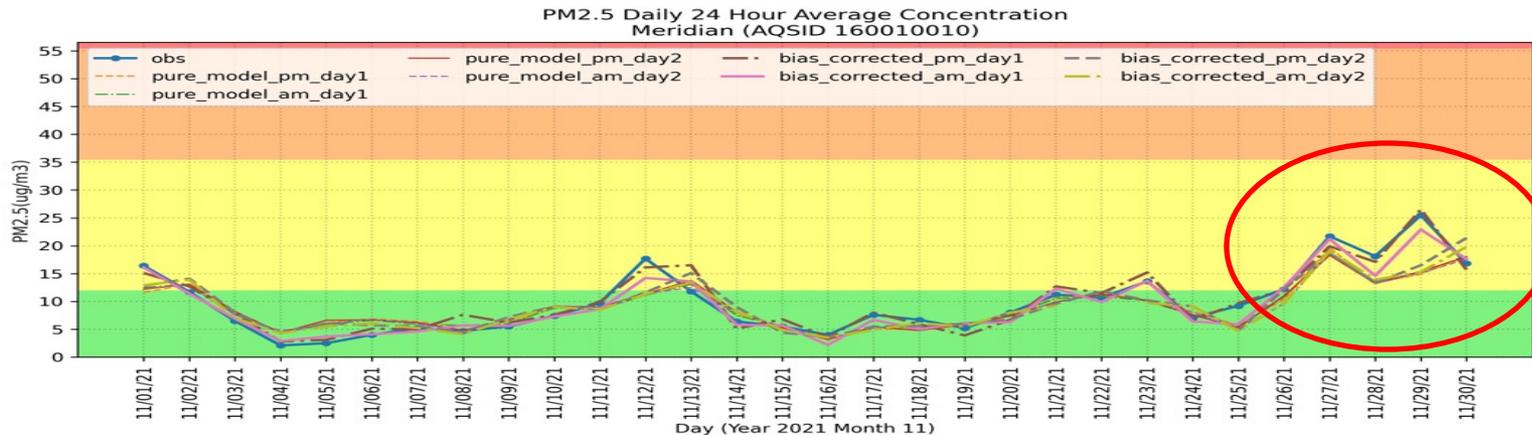
August 2021 Time Series



St. Lukes Meridian PM2.5 Site September 2021 Time Series

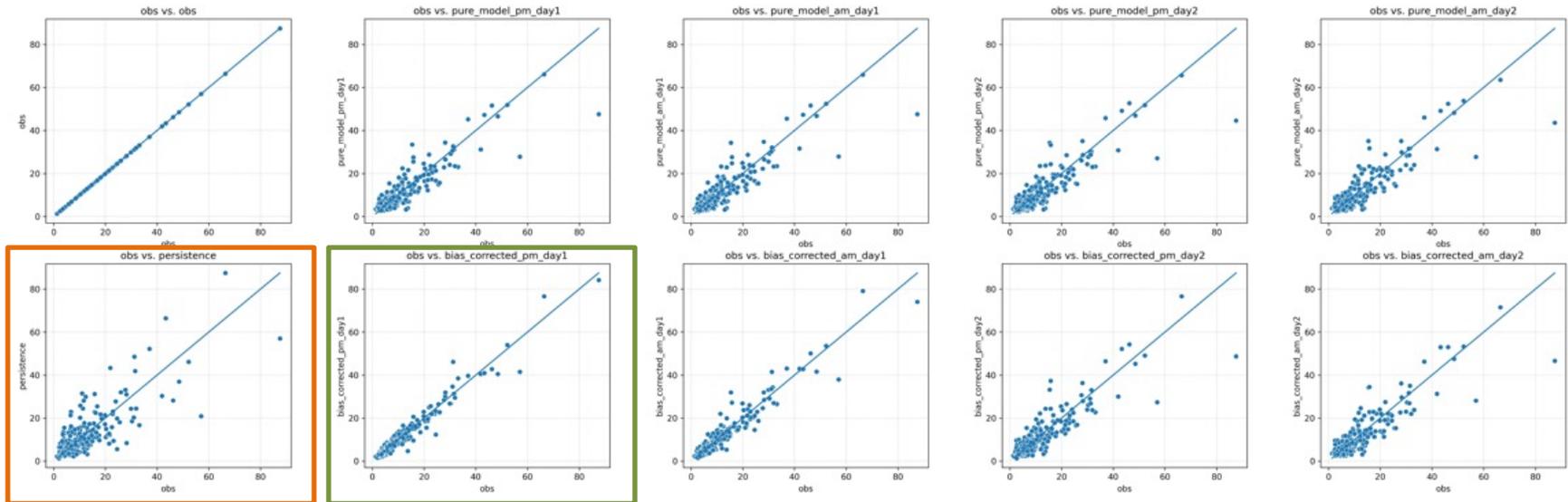


St. Lukes Meridian PM2.5 Site Nov. and Dec. 2021 Time Series



St. Lukes Meridian PM2.5 Site Year 2021 Daily Scatter Plots

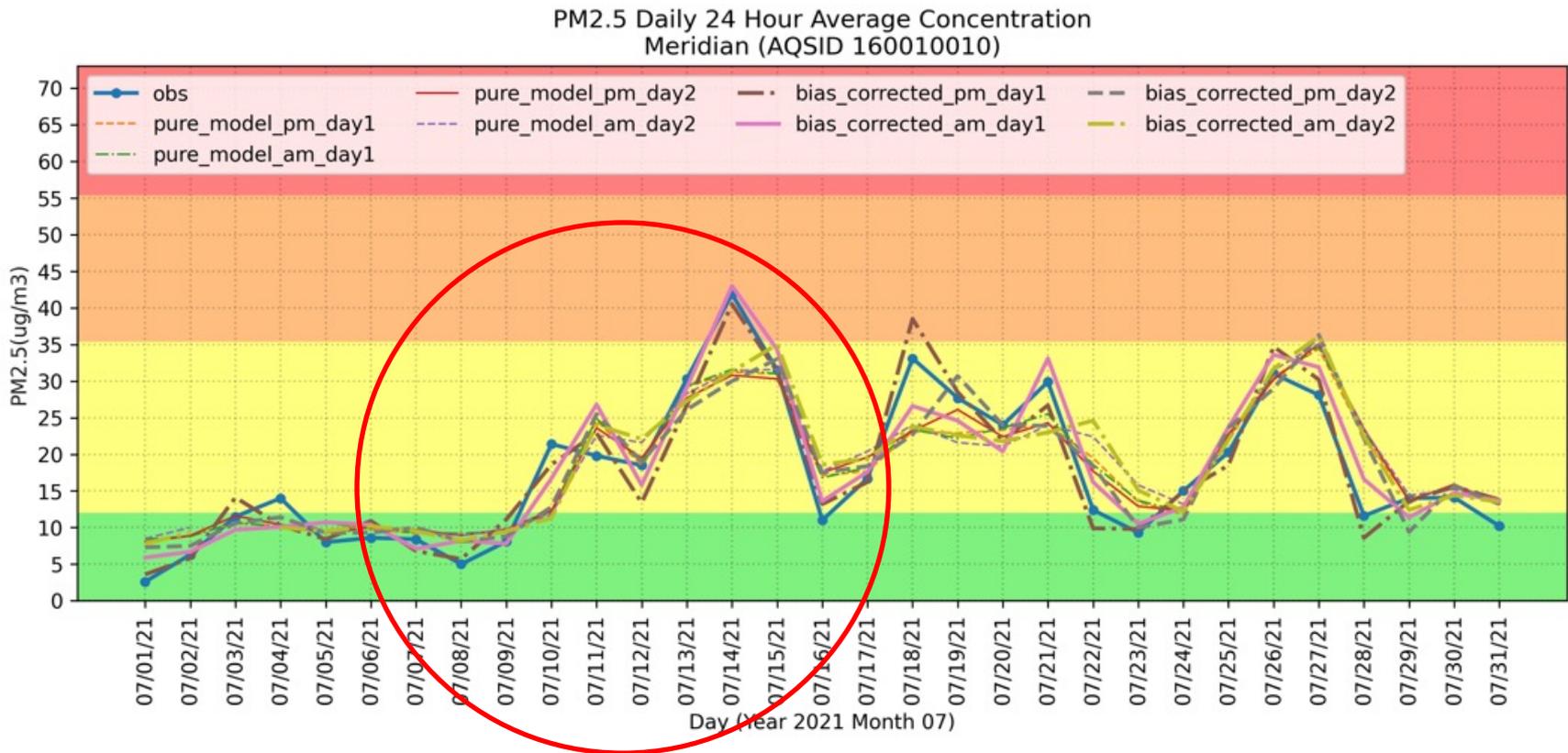
PM2.5 Daily 24 Hour Average Concentration Scatter Plots
Meridian (AQSID 160010010)
Year 2021



Wildfire Case Study
at
St. Lukes Meridian
for
PM2.5 Models

St. Lukes Meridian PM2.5 Site

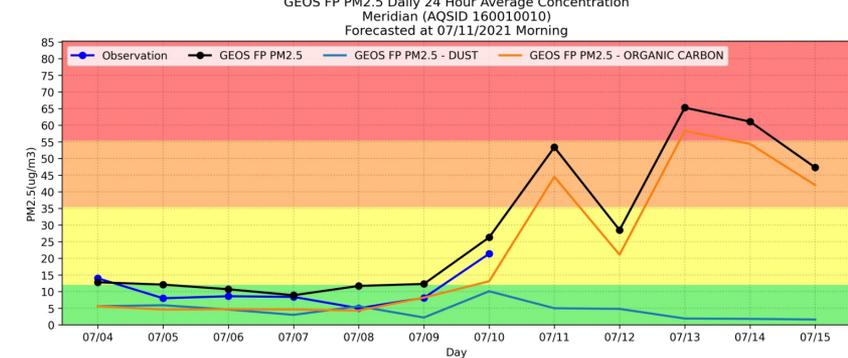
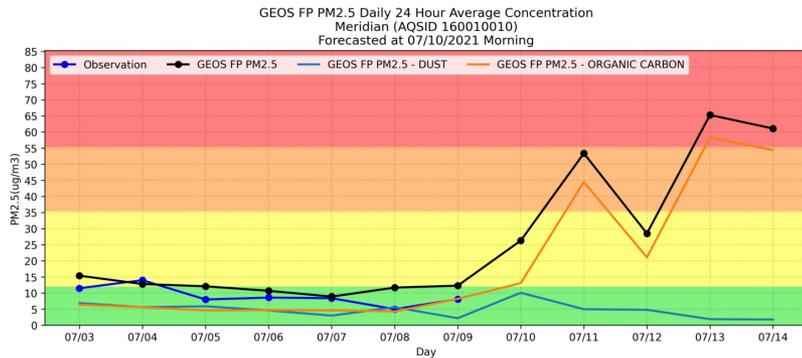
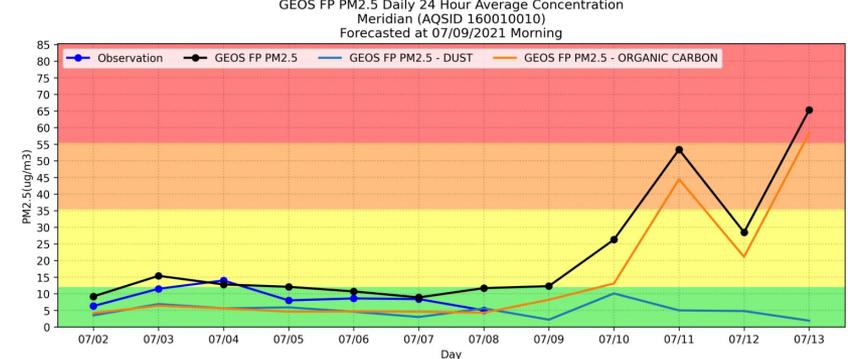
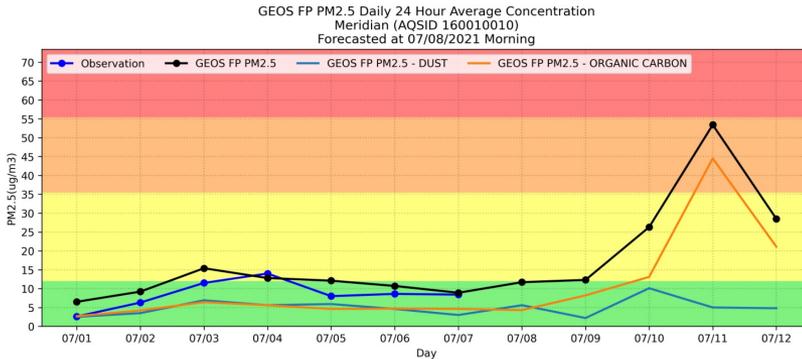
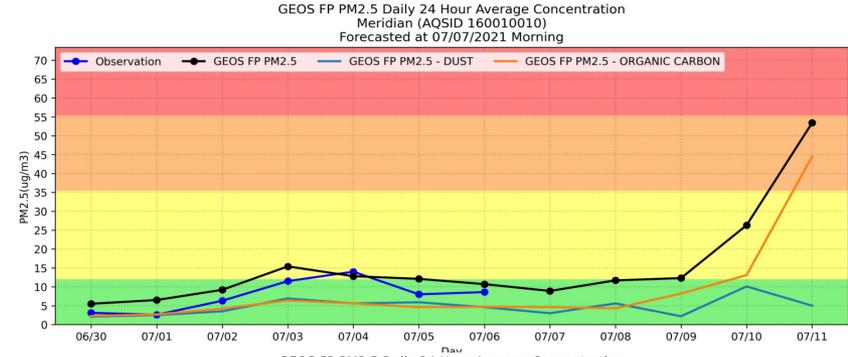
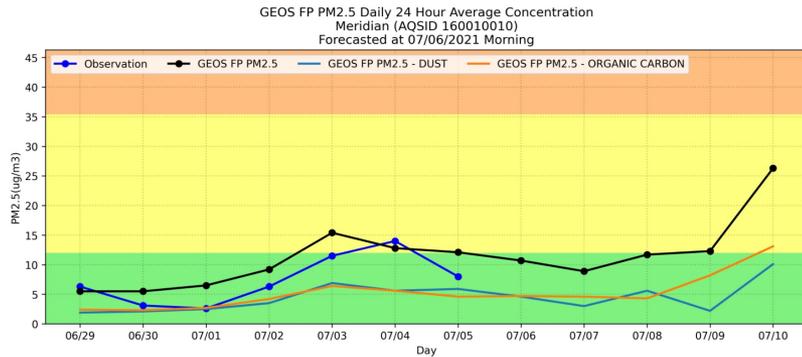
Wildfire Case Study Period : 07/6/2021 – 07/17/2021



NASA GEOS FP

Indication of Wildfire Impact

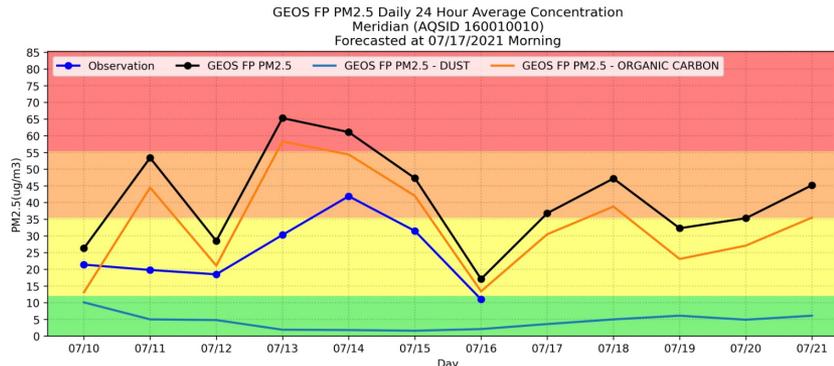
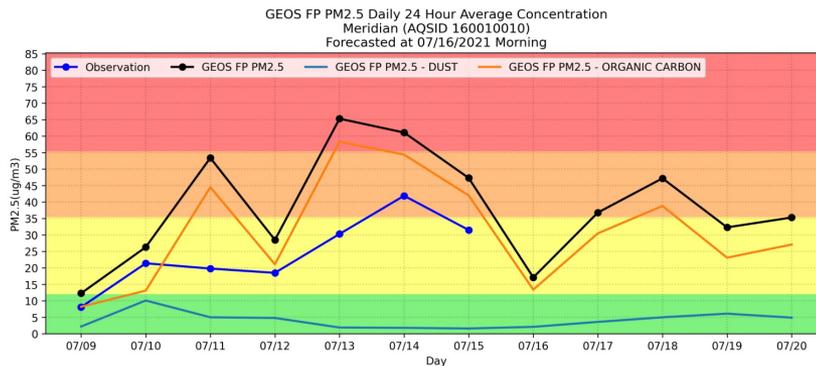
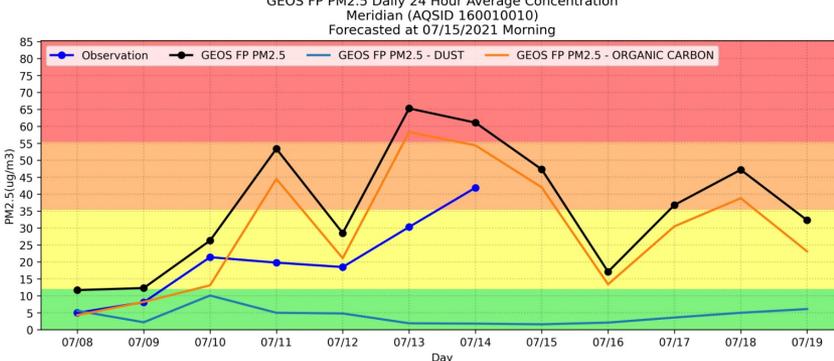
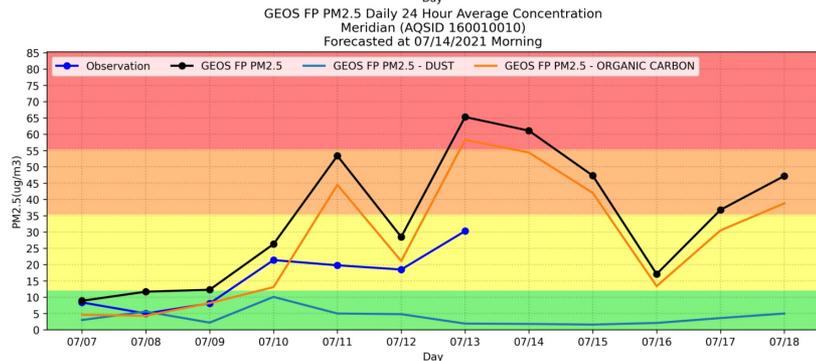
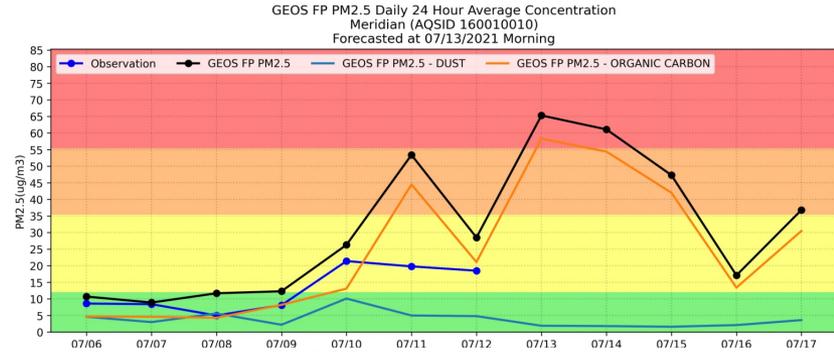
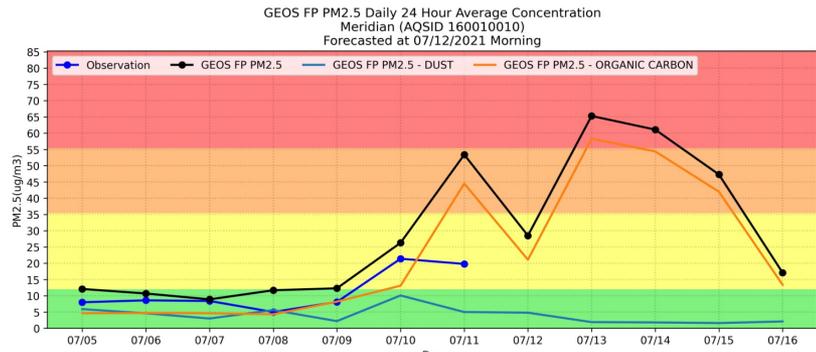
Morning Forecast on 7/6/2021 - 7/11/2021



NASA GEOS FP

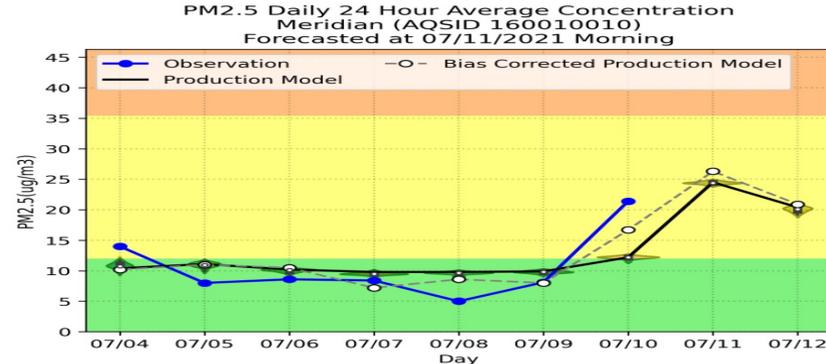
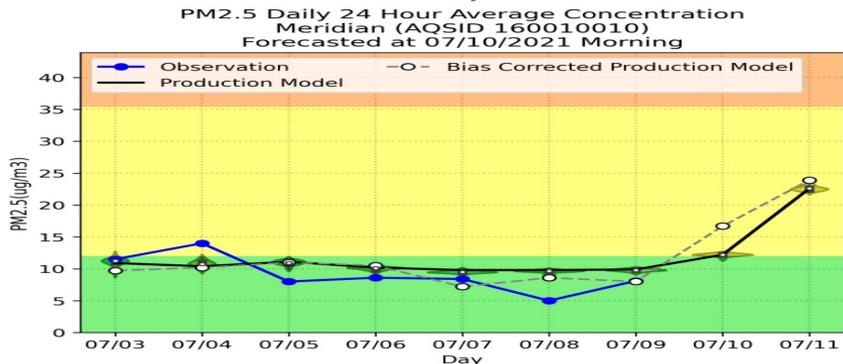
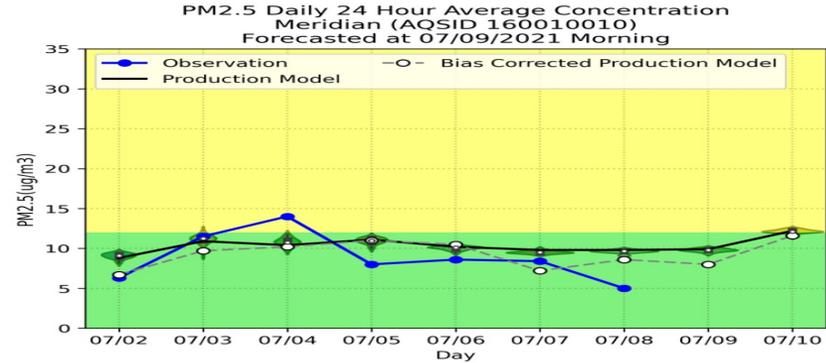
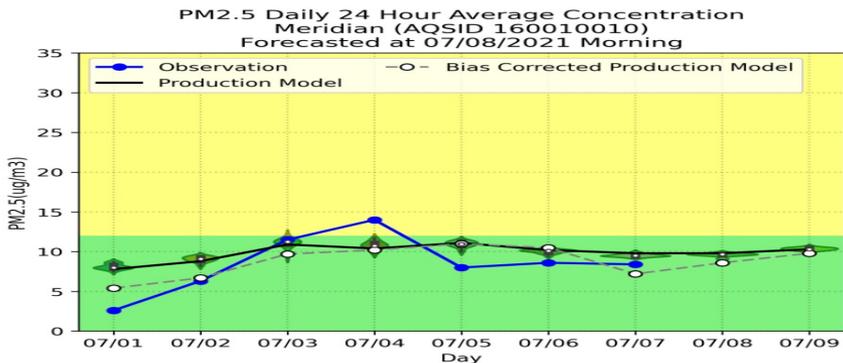
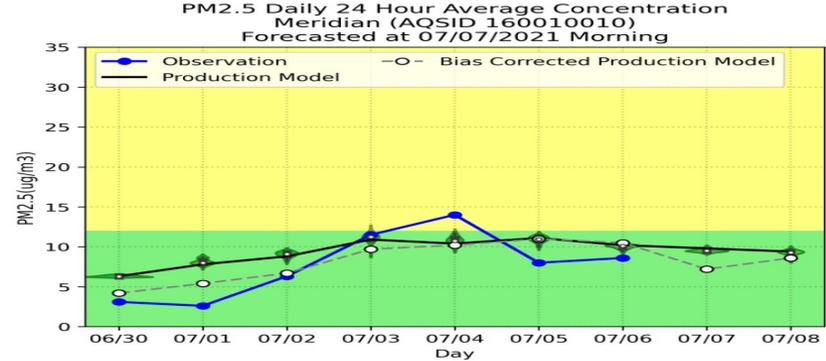
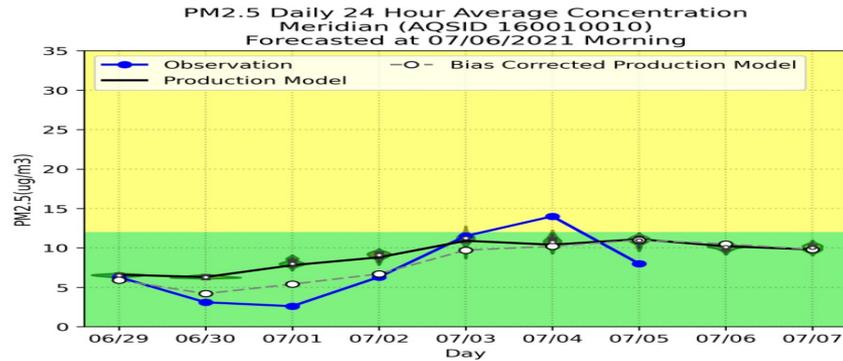
Indication of Wildfire Impact

Morning Forecast on 7/12/2021 - 7/17/2021



IDEQ Machine Learning Forecast System

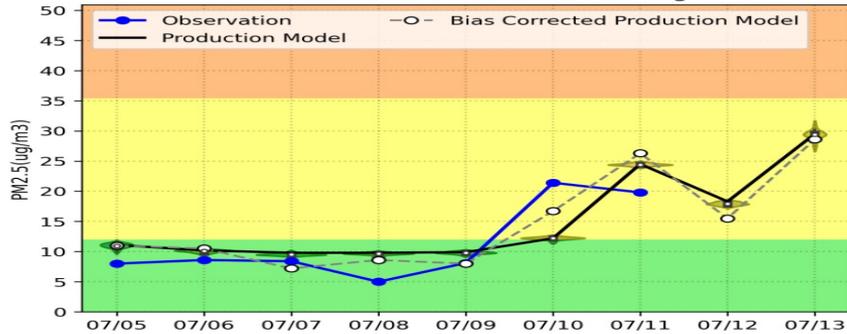
Morning Forecast on 7/6/2021 - 7/11/2021



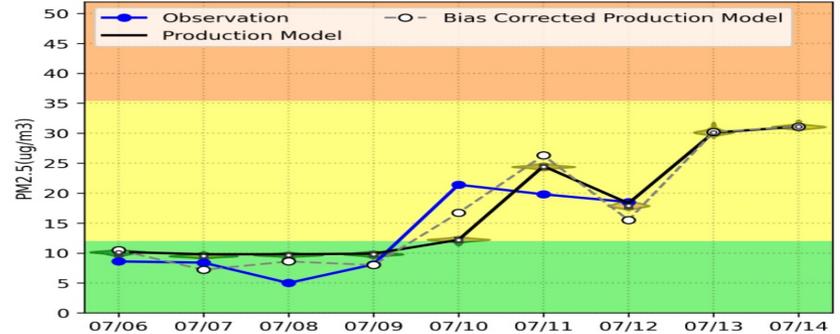
IDEQ Machine Learning Forecast System

Morning Forecast on 7/12/2021 - 7/17/2021

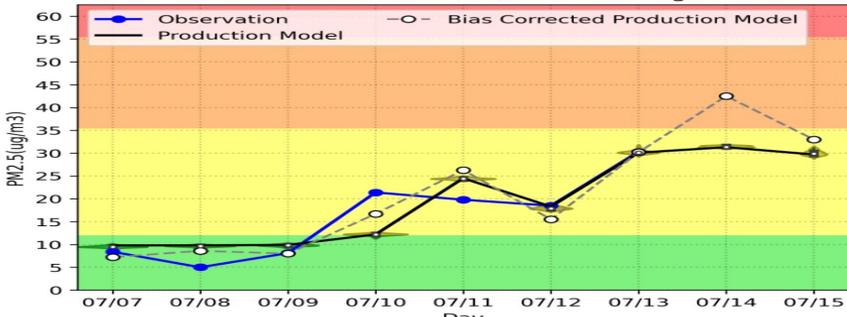
PM2.5 Daily 24 Hour Average Concentration
Meridian (AQSID 160010010)
Forecasted at 07/12/2021 Morning



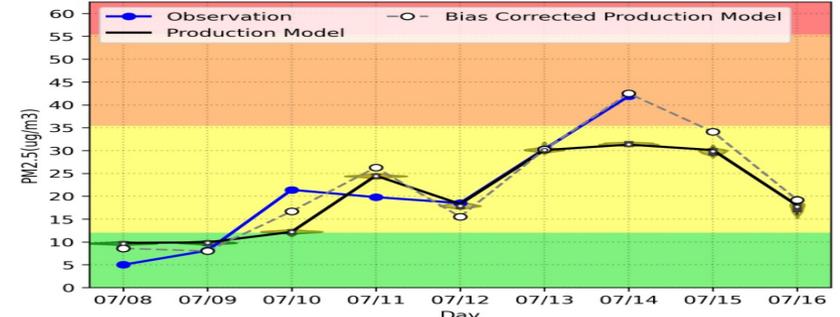
PM2.5 Daily 24 Hour Average Concentration
Meridian (AQSID 160010010)
Forecasted at 07/13/2021 Morning



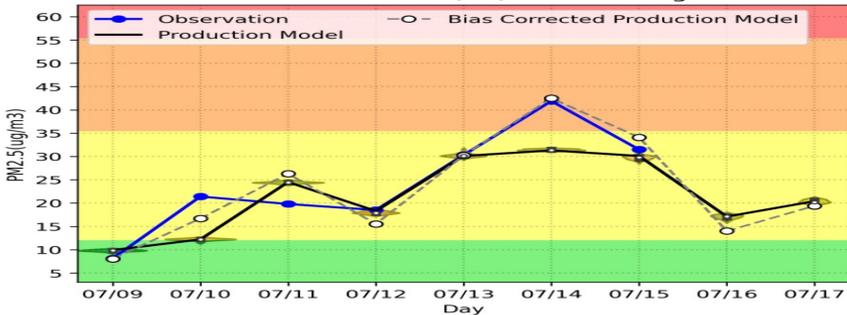
PM2.5 Daily 24 Hour Average Concentration
Meridian (AQSID 160010010)
Forecasted at 07/14/2021 Morning



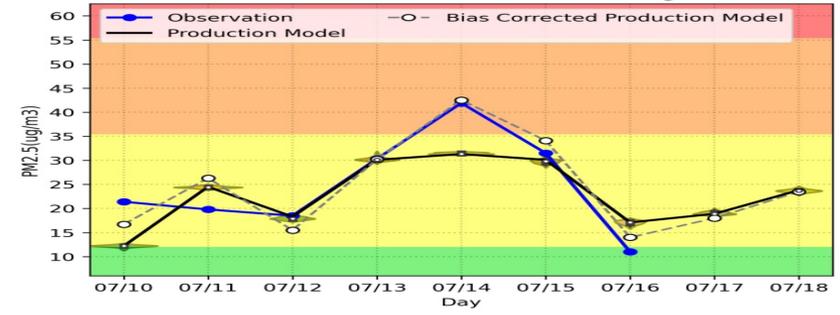
PM2.5 Daily 24 Hour Average Concentration
Meridian (AQSID 160010010)
Forecasted at 07/15/2021 Morning



PM2.5 Daily 24 Hour Average Concentration
Meridian (AQSID 160010010)
Forecasted at 07/16/2021 Morning

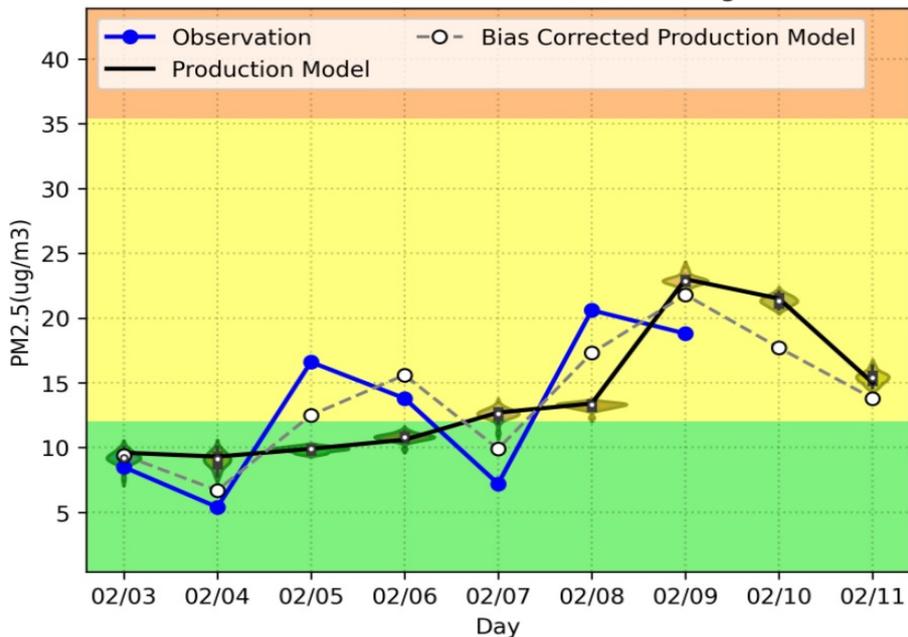


PM2.5 Daily 24 Hour Average Concentration
Meridian (AQSID 160010010)
Forecasted at 07/17/2021 Morning

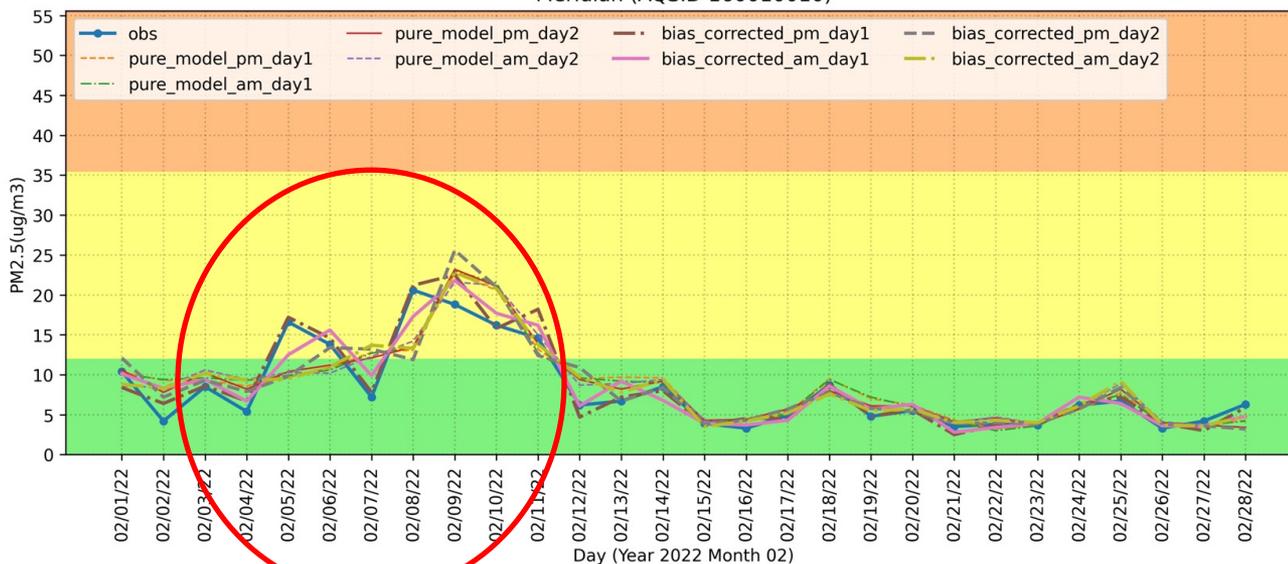


St. Lukes Meridian PM2.5 Site A Peek of Model Performance in Year 2022

PM2.5 Daily 24 Hour Average Concentration
Meridian (AQSID 160010010)
Forecasted at 02/10/2022 Morning



PM2.5 Daily 24 Hour Average Concentration
Meridian (AQSID 160010010)



Path Forward

- Annual Maintenance
 - Evaluate the model performance in January of each year
 - Retrain the machine learning models with last year's data added in January of each year
- Expand to forecast for day 3
- Explore the way to make prediction on unmonitored area
- Test and bring in new inputs to improve model performance in the future

Questions and Discussion

The End

Supplemental Slides

Email Notification

Machine Learning Air Quality Forecast is here! (Forecasted at 04/27/2022 Afternoon)



Wei.Zhang@deq.idaho.gov
To Wei Zhang



Wed 4/27/2022 1:48 PM

Wei,

This afternoon's forecast is here! The result can be found at:

http://10.220.98.54/ml_forecast_outputs/2022/20220427_pm/

OR through simple online map http://10.220.98.54/ml_forecast_outputs/2022/20220427_pm/20220427_pm_forecast_0Map.html

To access forecasts for previous days, please go to:

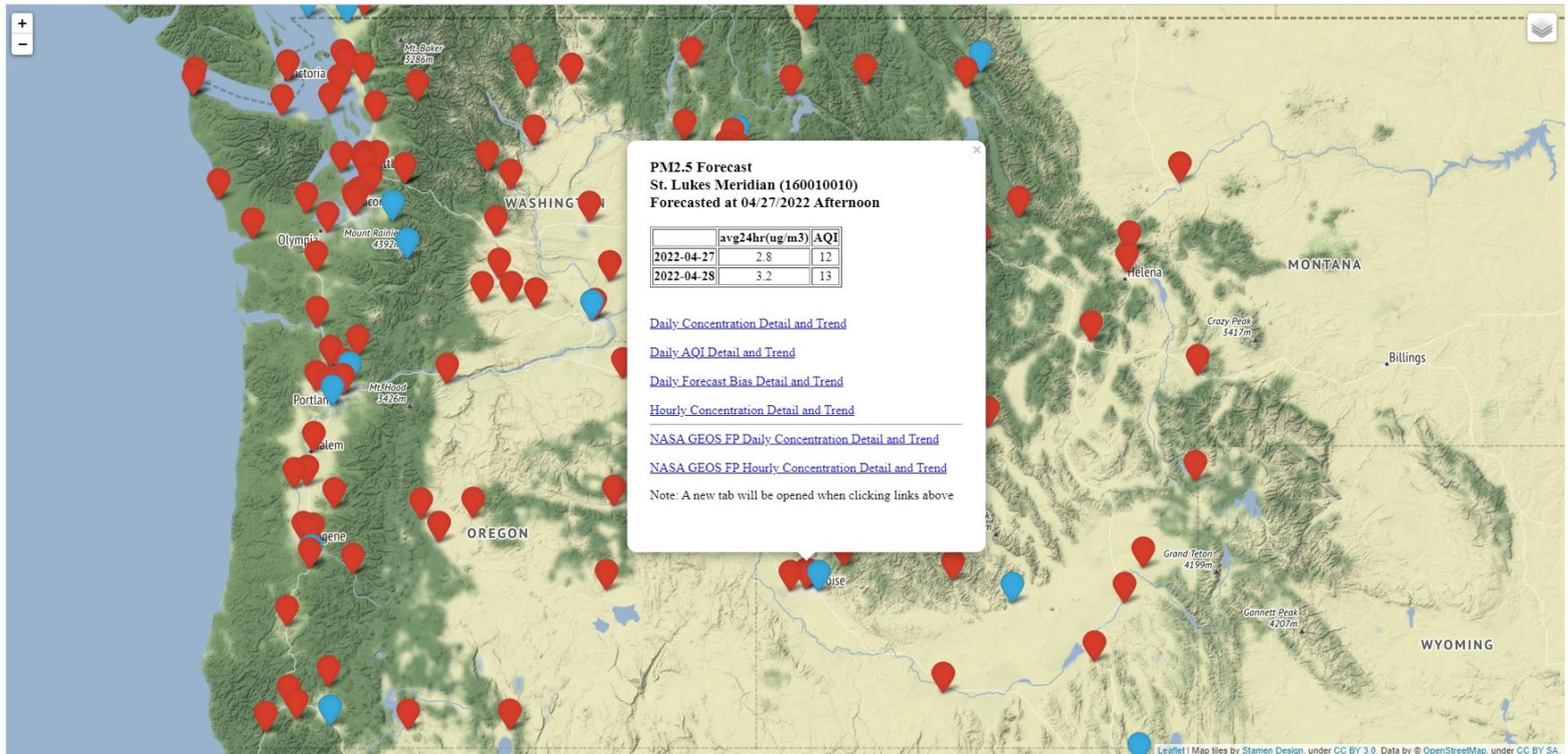
http://10.220.98.54/ml_forecast_outputs/

This forecast is based on a machine learning method utilizing modeled ensemble meteorological forecasts from the University of Washington. Acute or event-based impacts, such as wildfires and dust storms, are not considered in the model. This forecast should be used as a starting point and then adjusted based on local knowledge of these factors.

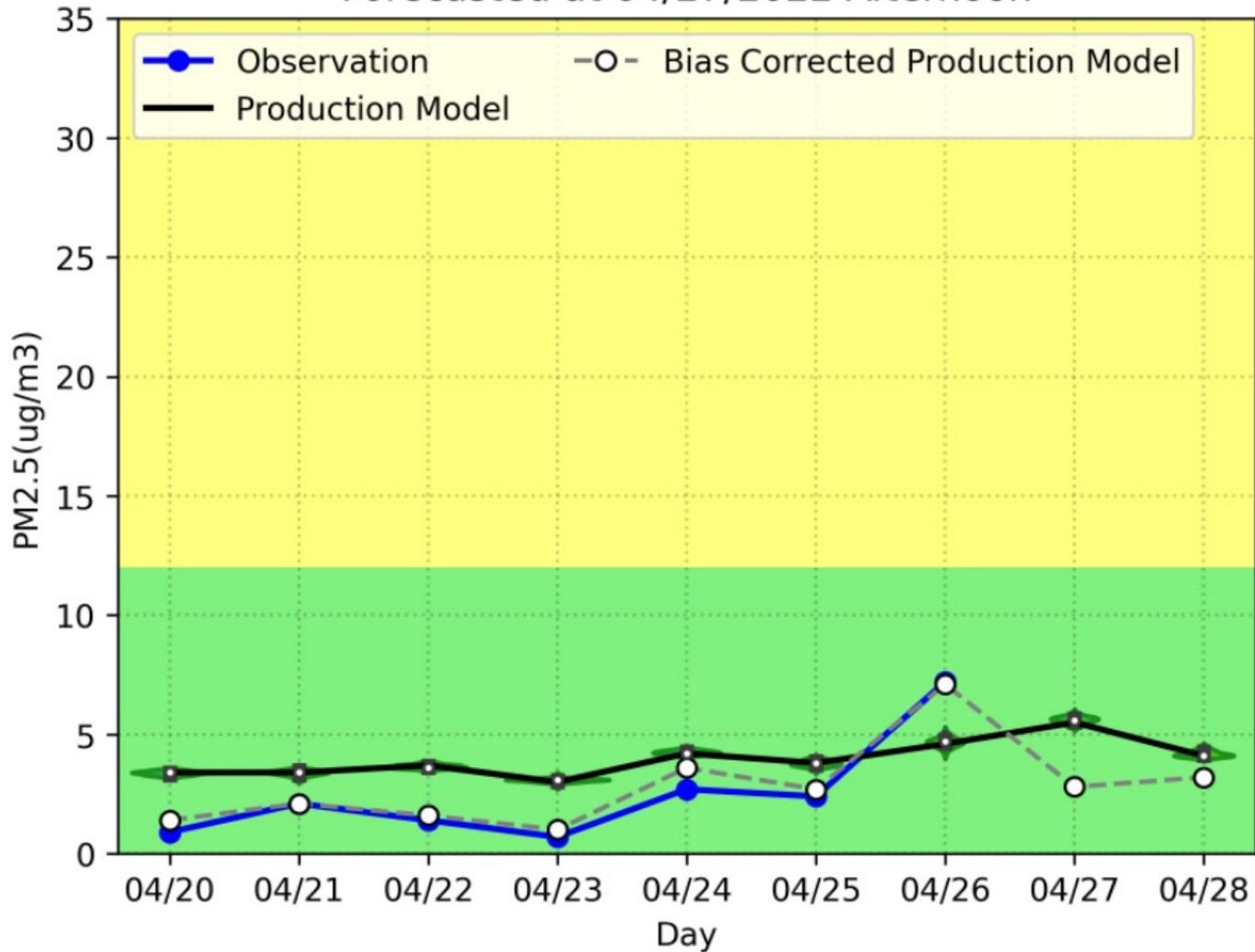
Hopefully this tool makes your life easier. Please provide us any feedback you may have.

Technical Services Division
Idaho Department of Environmental Quality
Wei.Zhang@deq.idaho.gov

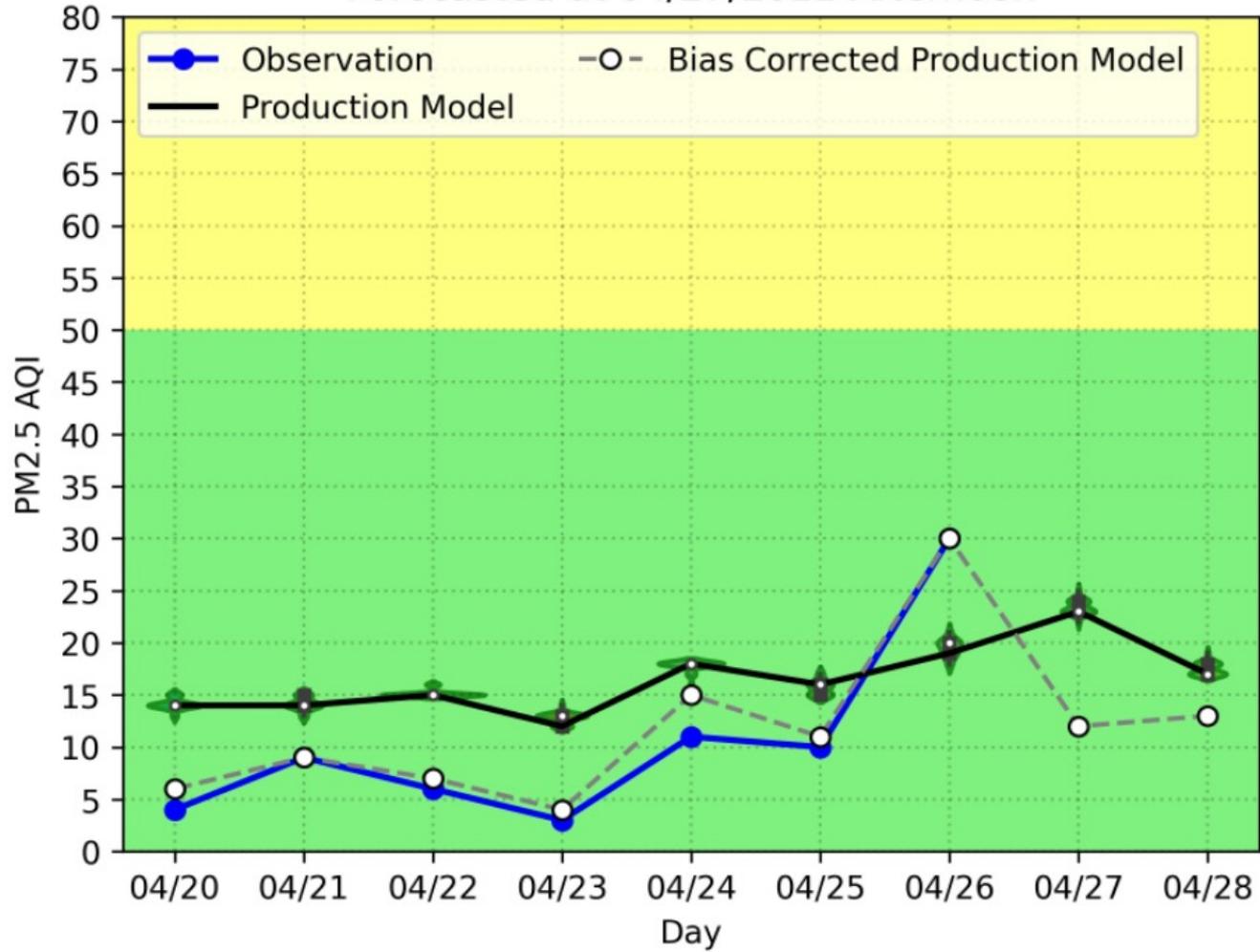
Simple Online Map



PM2.5 Daily 24 Hour Average Concentration
Meridian (AQSID 160010010)
Forecasted at 04/27/2022 Afternoon

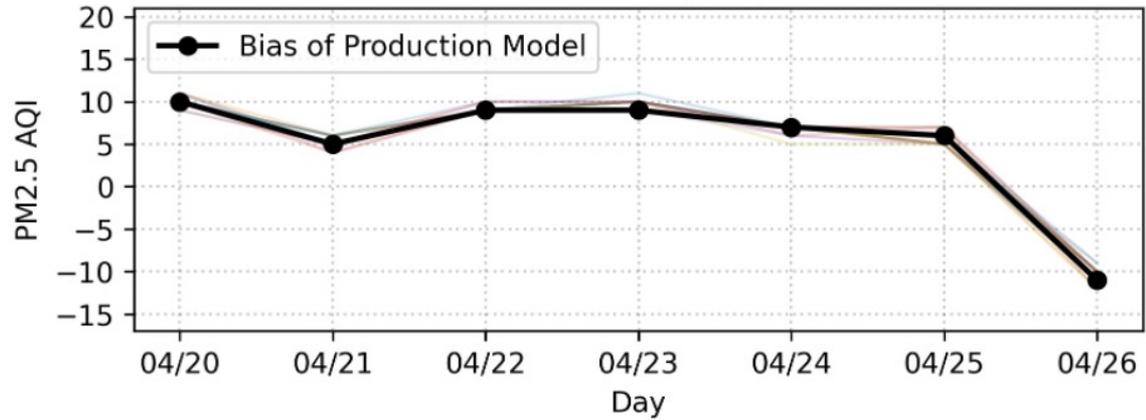
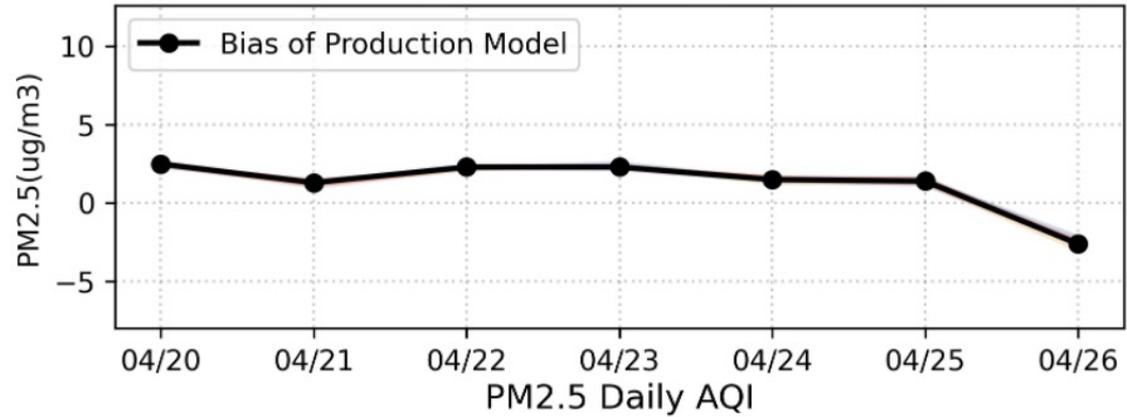


PM2.5 Daily AQI
Meridian (AQSID 160010010)
Forecasted at 04/27/2022 Afternoon

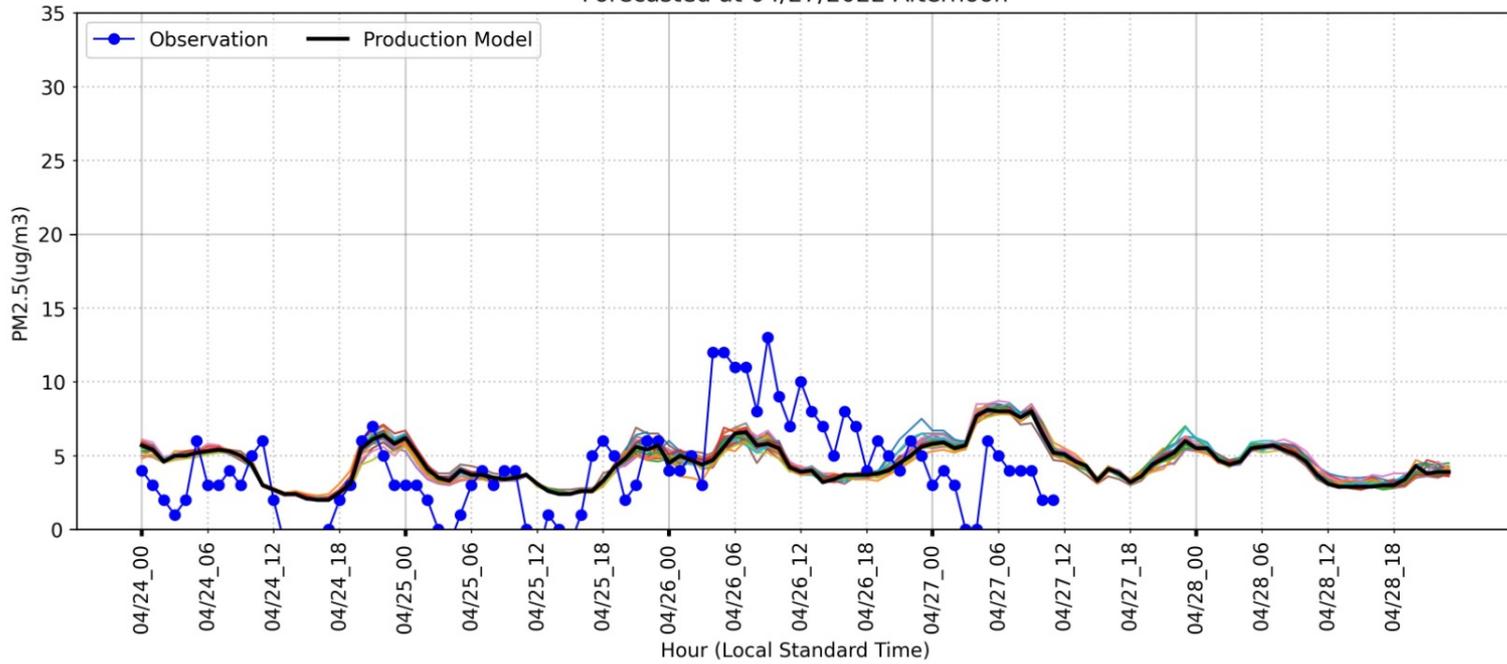


PM2.5 Daily Forecast Bias
Meridian (AQSID 160010010)

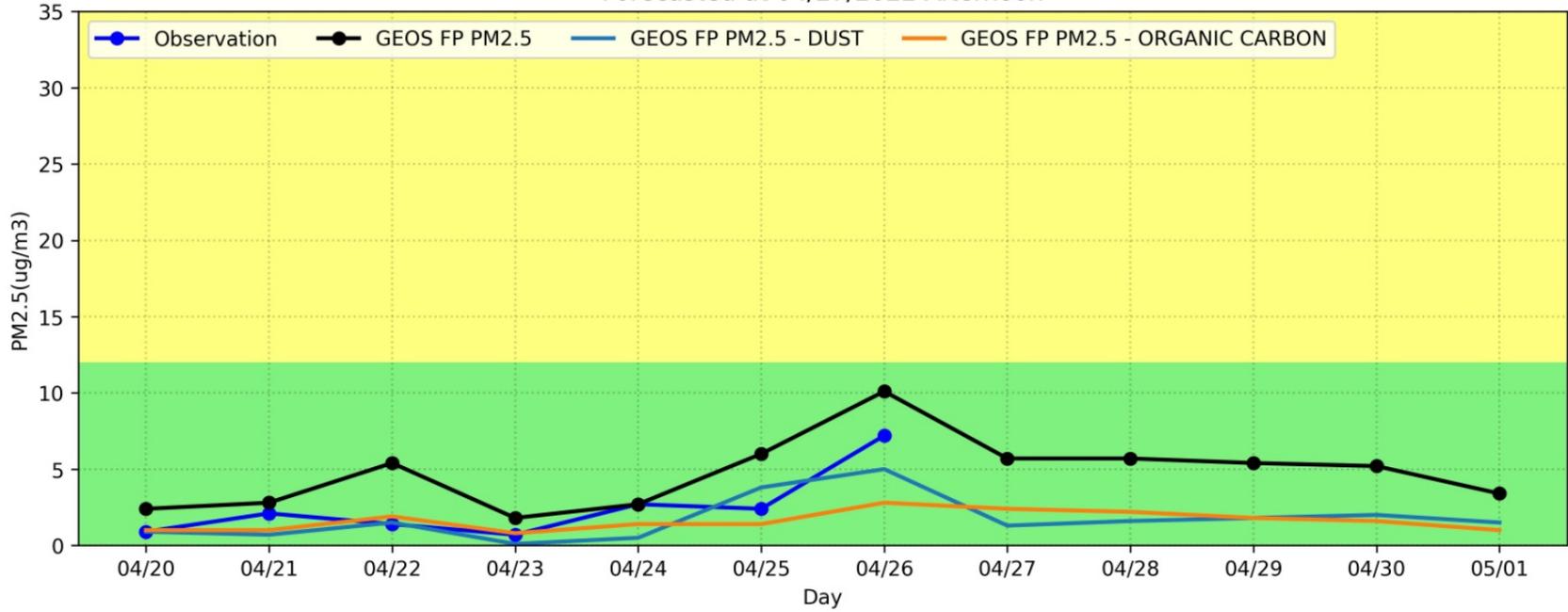
PM2.5 Daily 24 Hour Average Concentration



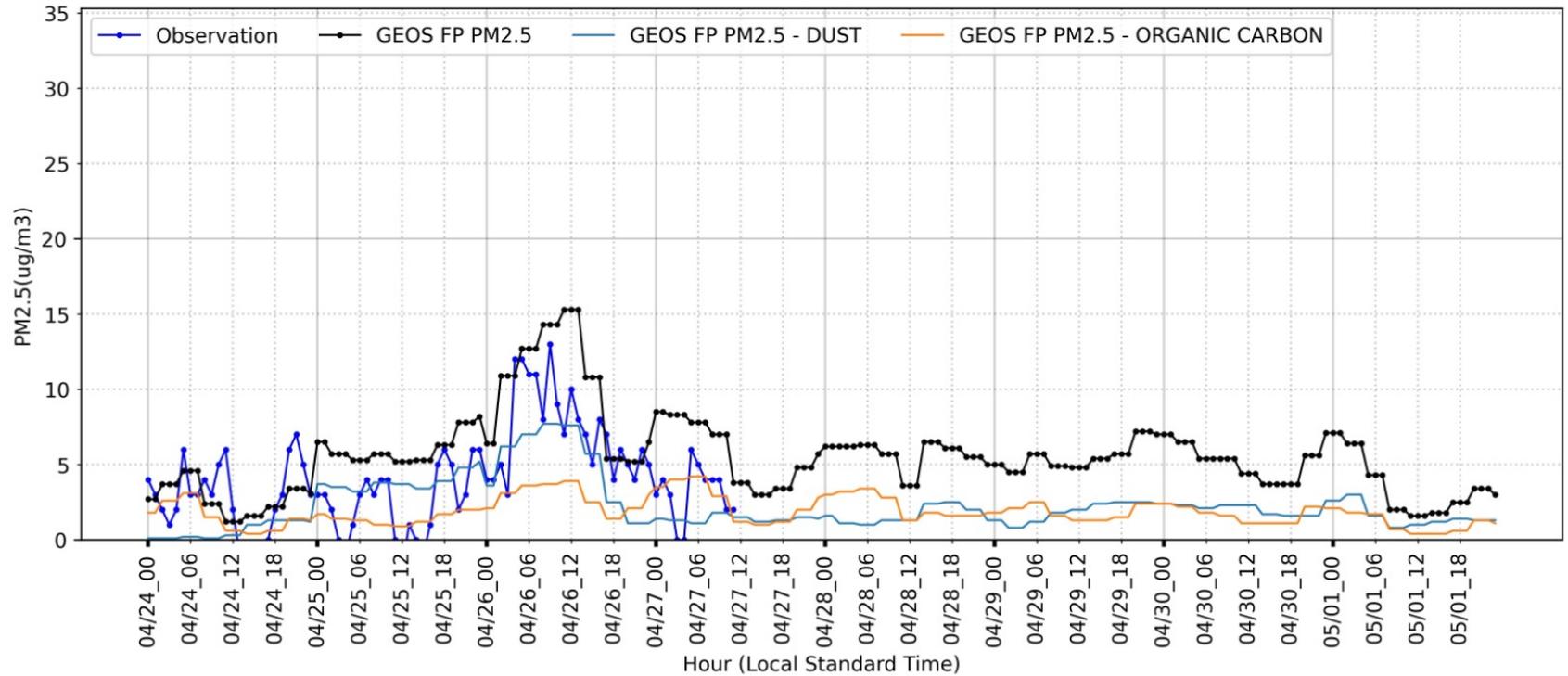
PM2.5 Hourly Concentration
Meridian (AQSID 160010010)
Forecasted at 04/27/2022 Afternoon



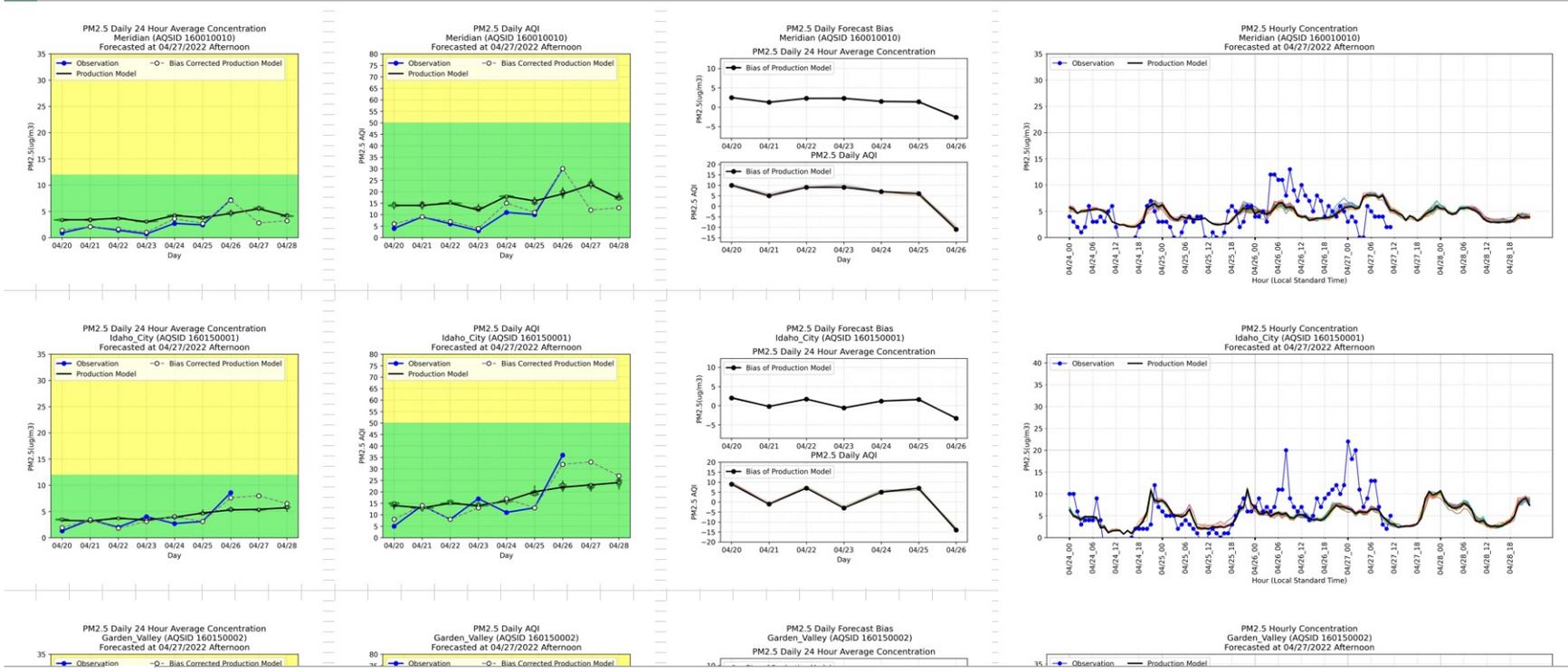
GEOS FP PM2.5 Daily 24 Hour Average Concentration
Meridian (AQSID 160010010)
Forecasted at 04/27/2022 Afternoon



GEOS FP PM2.5 Hourly Concentration
Meridian (AQSID 160010010)
Forecasted at 04/27/2022 Afternoon

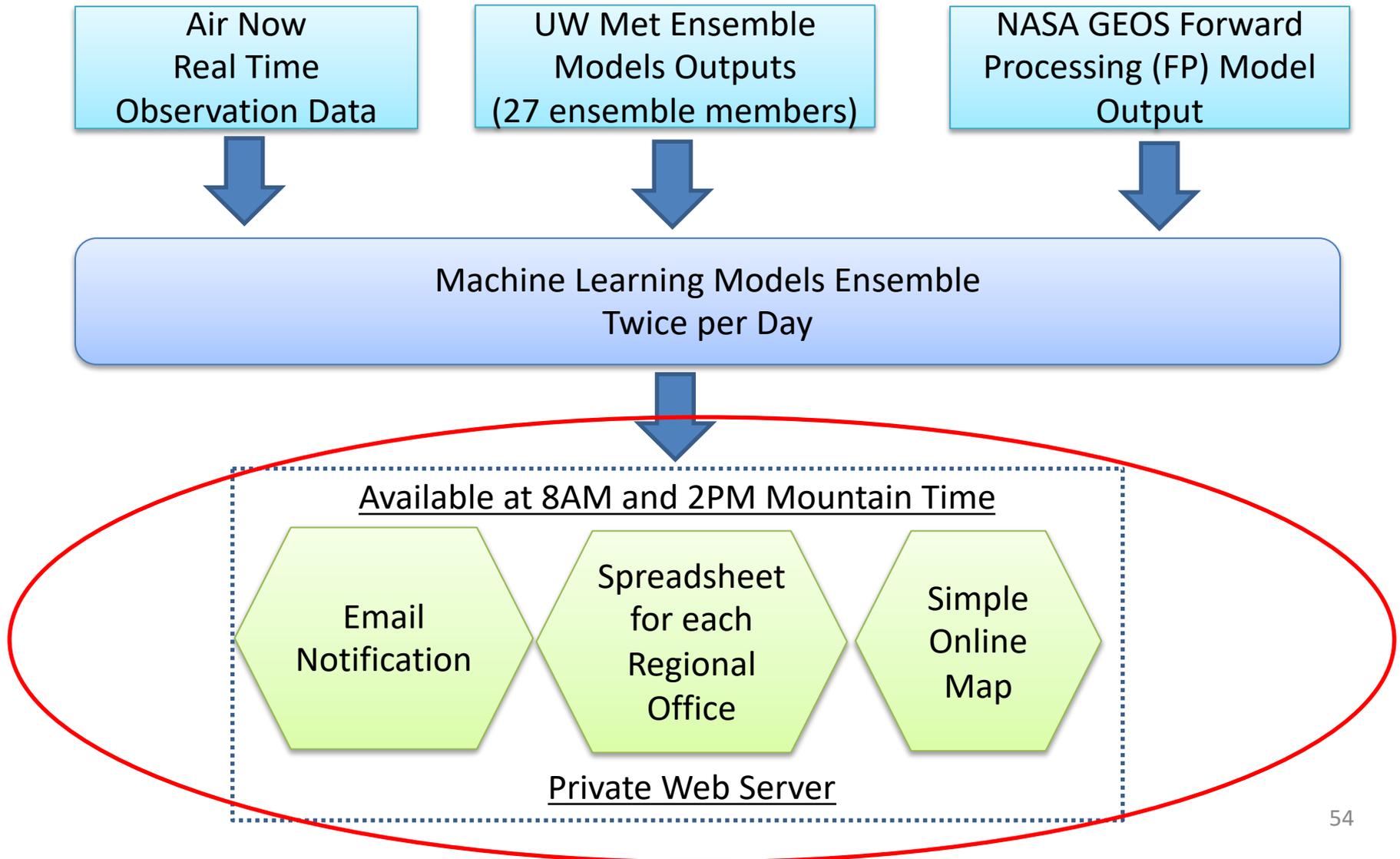


Spreadsheet per Regional Office



Overview of Forecast System

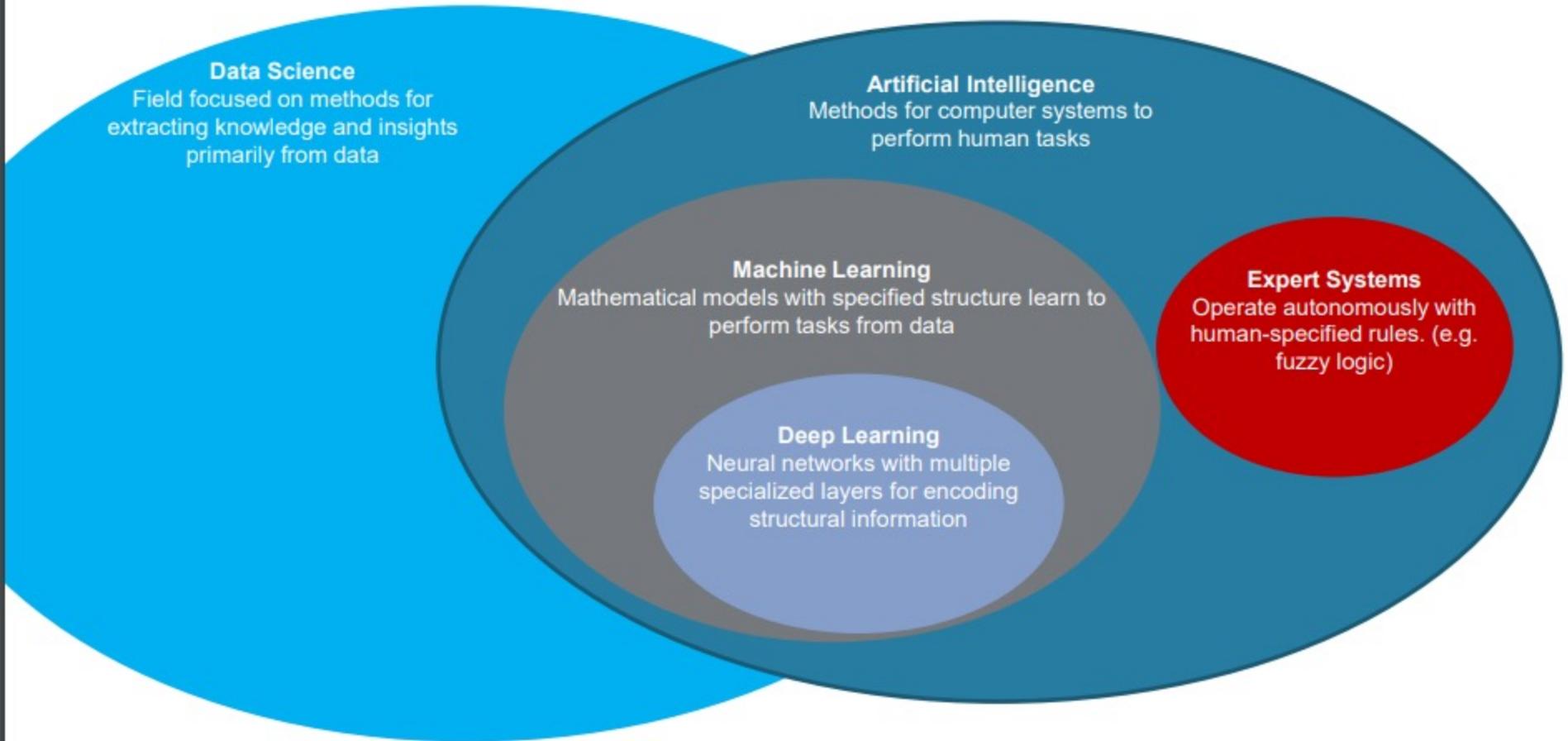
Data Process – In and Out



Demo

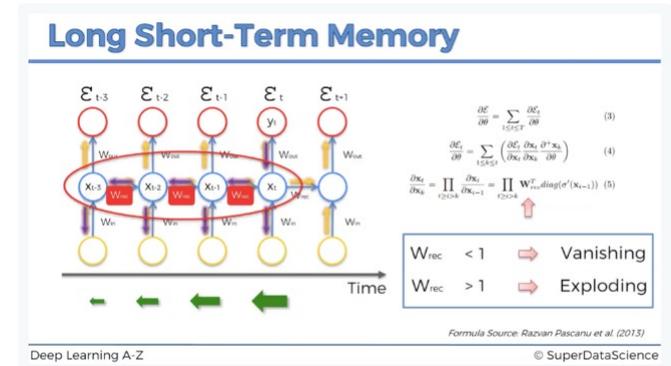
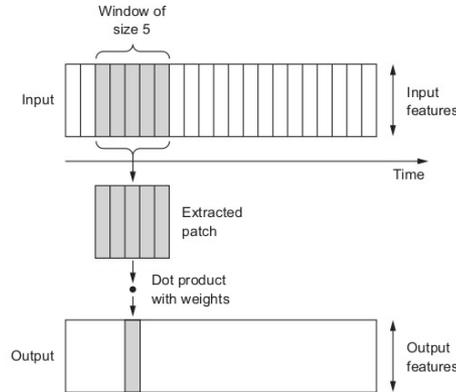
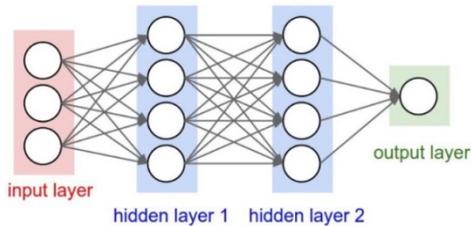
- Email Notification
- Simple Online Map
 - Make it easy to navigate the sites
- Spreadsheet for each Regional Office
 - Regional offices pick the sites

The Data Science Taxonomy



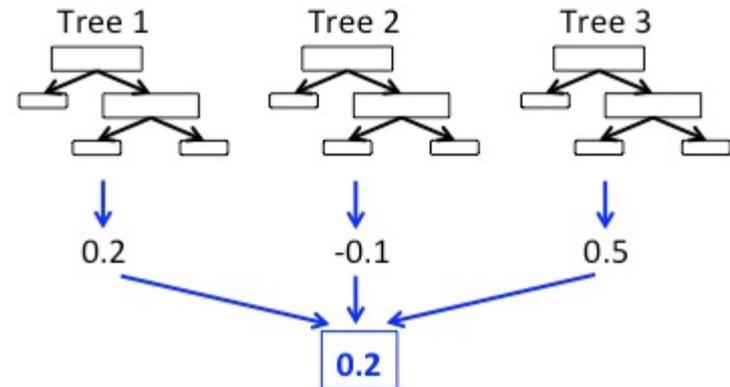
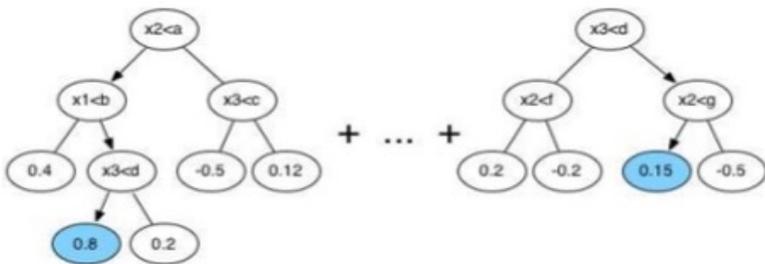
Neural Network

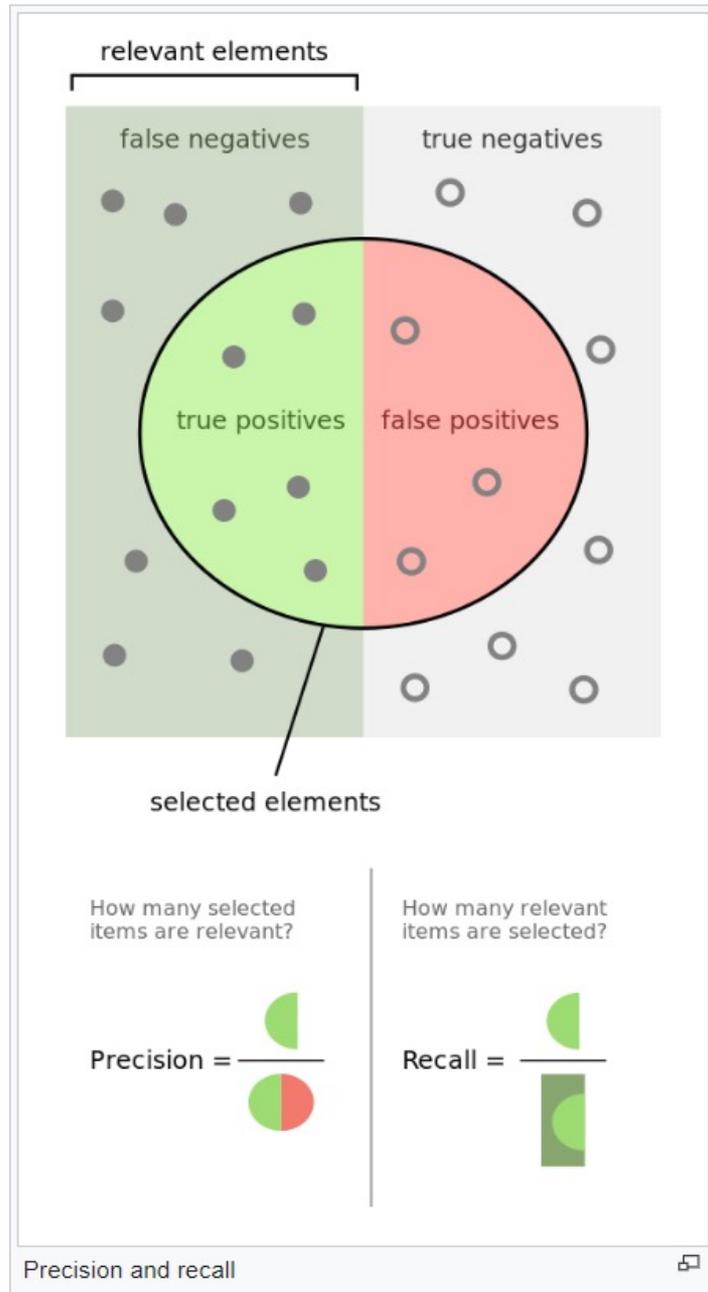
- Dense Neural Network
- 1D Convolutional Neural Network
- Recurrent Neural Network (LSTM)



Tree based Methods

- XGBoost
 - XGBoost stands for e**X**treme **G**radient **B**oosting
 - Tree built sequentially by minimizing the residue (error) of the previous tree
- Random Forest
- Boosted Random Forest





Precision Recall F1 Score

$$F = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

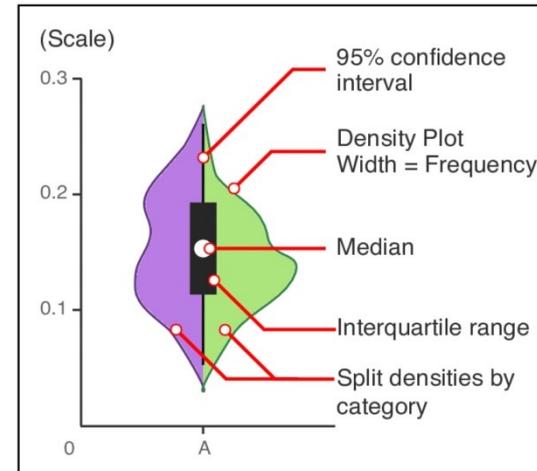
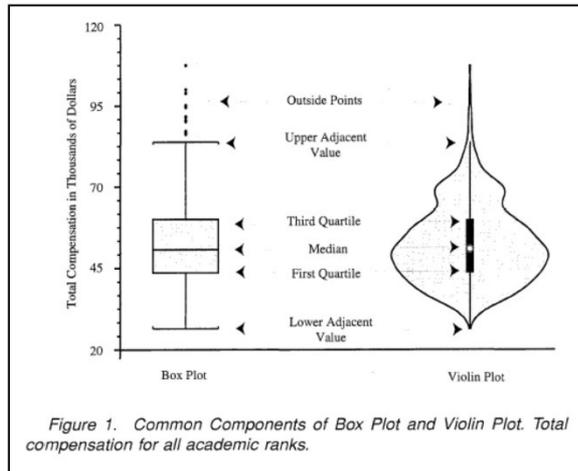
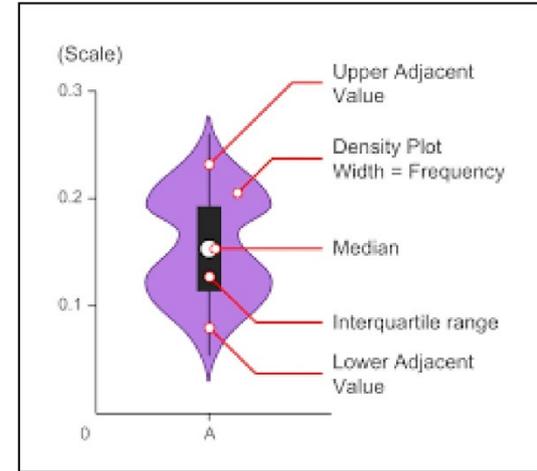
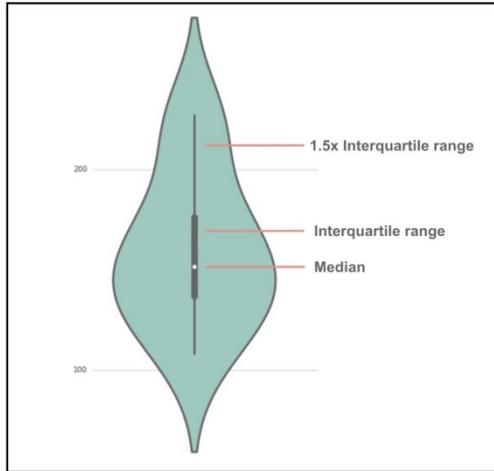
A measure that combines precision and recall is the [harmonic mean](#) of precision and recall, the traditional F-measure or balanced F-score

Heidke Skill Score (HSS) and Hanssen-Kuiper Skill Score (KSS)

- HSS represents the accuracy of the model prediction compared with a reference forecast, which is from the random guess that is statistically independent of the observations.
- The range of the HSS is from $-\infty$ to 1. A negative value means a random guess is better, 0 means no skill, and 1 means a perfect score.
- KSS measures the ability to separate different categories. The range is from -1 to 1 where 0 means no skill, and 1 means a perfect score.

Violin Plot Explained

Present Distribution of Model Ensemble Member Forecasts



NASA GEOS

GEOS Forecast and Reanalysis Products

GEOS Forward Processing
(GEOS FP)
NRT Analysis and Forecast

GEOS-Composition
Forecast
(GEOS-CF)
NRT Forecast

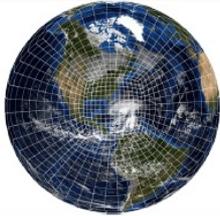
Modern-Era Retrospective
analysis for Research and
Applications, Version 2
(MERRA-2) Reanalysis



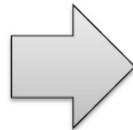
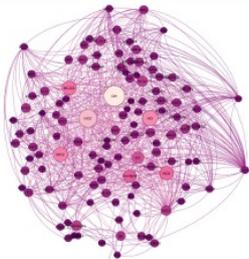
NASA GEOS-CF

NASA Composition Forecasts (GEOS-CF)

GEOS
Meteorology



GEOS-Chem
Chemistry



GEOS-CF

- [GEOS-Chem](#): Global chemistry transport model driven by GEOS meteorology
- 1-day simulation of the previous day using the analysis from FP-IT
 - Uses a **replay** technique to force the meteorology towards the FP-IT analysis
 - FP-IT is a 'frozen' version of FP used for satellite retrievals, similar to the version used to make MERRA-2.
- 5-day forecast
- Two aerosol schemes:
 - GOCART – Radiatively coupled to AGCM
 - GEOS-Chem – No feedbacks to model physics
- Full description in [Keller et al., 2021](#)

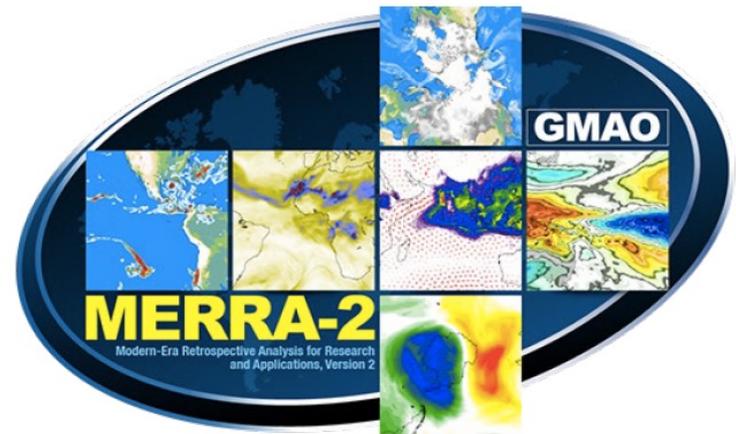


NASA GEOS MERRA-2

MERRA-2 Reanalysis

<https://gmao.gsfc.nasa.gov/reanalysis/MERRA-2/>

- The **M**odern-**E**ra **R**etrospective analysis for **R**esearch and **A**pplications version **2** (MERRA-2) provides data beginning in 1980 and runs a few weeks behind real-time.
- Long-term, model-based analyses of multiple datasets using a fixed assimilation system
- Includes meteorology, stratospheric ozone, and aerosols at the spatial resolution of a $0.5^\circ \times 0.66^\circ$ (~50 km) grid.



Source: <https://gmao.gsfc.nasa.gov/reanalysis/>



NASA Goddard Earth Observing System (GEOS) modeling and data assimilation systems

GEOS Output Quick Guide

	GEOS FP	GEOS-CF	MERRA-2
Type	Analysis + Forecast	Replay + Forecast	Reanalysis
Domain	Global	Global	Global
Spatial Resolution	Simulation: ~12 km Output: ~25 km (0.25°x0.312°)	~25 km (0.25°x0.312°)	~50km (0.5°x0.625°)
Temporal Resolution	2-D data: Hourly 3-D data: Every 3 h	15 min, Hourly	Hourly, Daily, Monthly
Vertical Levels	72 (near surface-0.1 hPa)	72 (near surface-0.1 hPa)	72 (near surface-0.1 hPa)
Output available	Analysis: 2014 – Present Forecast: ~20 days	Replay: 2018 – Present Forecast: 2019 – Present (aqc collection) ~14 days (all collections)	1980-Present
Initialization	Daily 10-day forecast at 00Z Daily 5-day forecast at 12Z	Daily 5-day forecast at 12Z	~1-2 months behind real time
Data Assimilation	Yes	No	Yes
File Specification Doc	https://gmao.gsfc.nasa.gov/pubs/docs/Lucchesi1203.pdf *	https://gmao.gsfc.nasa.gov/pubs/docs/Knowland1204.pdf *	https://gmao.gsfc.nasa.gov/pubs/docs/Bosilovich785.pdf *

NASA's Applied Remote Sensing Training Program

* Find most current File Specification at https://gmao.gsfc.nasa.gov/pubs/office_notes.php

43

[ARSET - Introduction and Access to Global Air Quality Forecasting Data and Tools | NASA Applied Science](#)

Terrain Feature

- Source : Esri
- World Landforms - Improved Hammond Method
- 16 classes of landform types and regions
 - Nearly flat plains
 - Smooth plains with some local relief
 - Irregular plains with moderate relief
 - Irregular plains with low hills
 - Scattered moderate hills
 - Scattered high hills
 - Scattered low mountains
 - Scattered high mountains
 - Moderate hills
 - High hills
 - Tablelands with moderate relief
 - Tablelands with considerable relief
 - Tablelands with high relief
 - Tablelands with very high relief
 - Low mountains
 - High mountains