

IDEQ Machine Learning Air Quality Forecast System Introduction

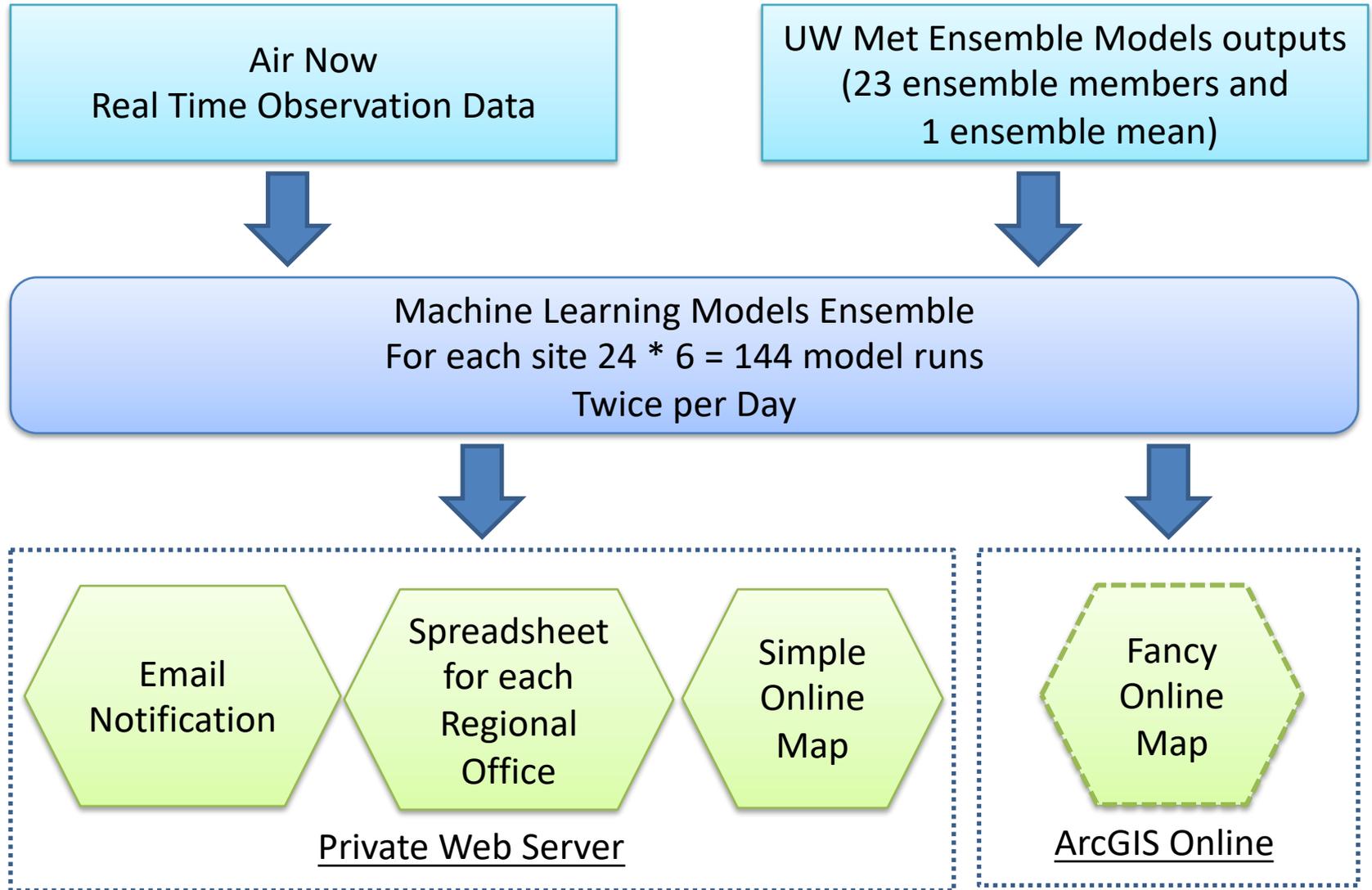
NW-AIRQUEST 2021 Annual Meeting

06/08/2021

Acknowledgment

- Kai Fan from Washington State University (WSU) and Ranil Dhammapala from Washington Department of Ecology (WECY) provided us the 2017-2019 data set and shared their experiences with us
- David Ovens from University of Washington (UW) provided the real time Meteorology Ensemble Models outputs
- Air Now provided the real time observation data
- NRMC and NW-AIRQUEST

Overview of Forecast System



UW Meteorology Models

- Currently 23 ensemble members and 1 ensemble mean
- 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

O3 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

- Pure XGBoost
- Random Forest
- Boosted Random Forest

Baseline Model : Persistence

O3 Forecast

Data Used in Model Development

- From Kai and Ranil (Thanks!)
- Meteorology Parameters :
 - UW WRF Production Run
- Training Data Set : 2017 and 2018
- Validation : K-Fold
- Testing Data Set : 2019

O3 Forecast Features Used

11 Features :

- Previous 24th hour O3
- T, P, RH, PBLH, WS, WDIR
- O3hourly_mean (grouped by Month, Weekday/end, Hour)
- Month (1-12), Weekday (0-1), Hour (0-23)

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

St. Lukes Meridian O3 Site

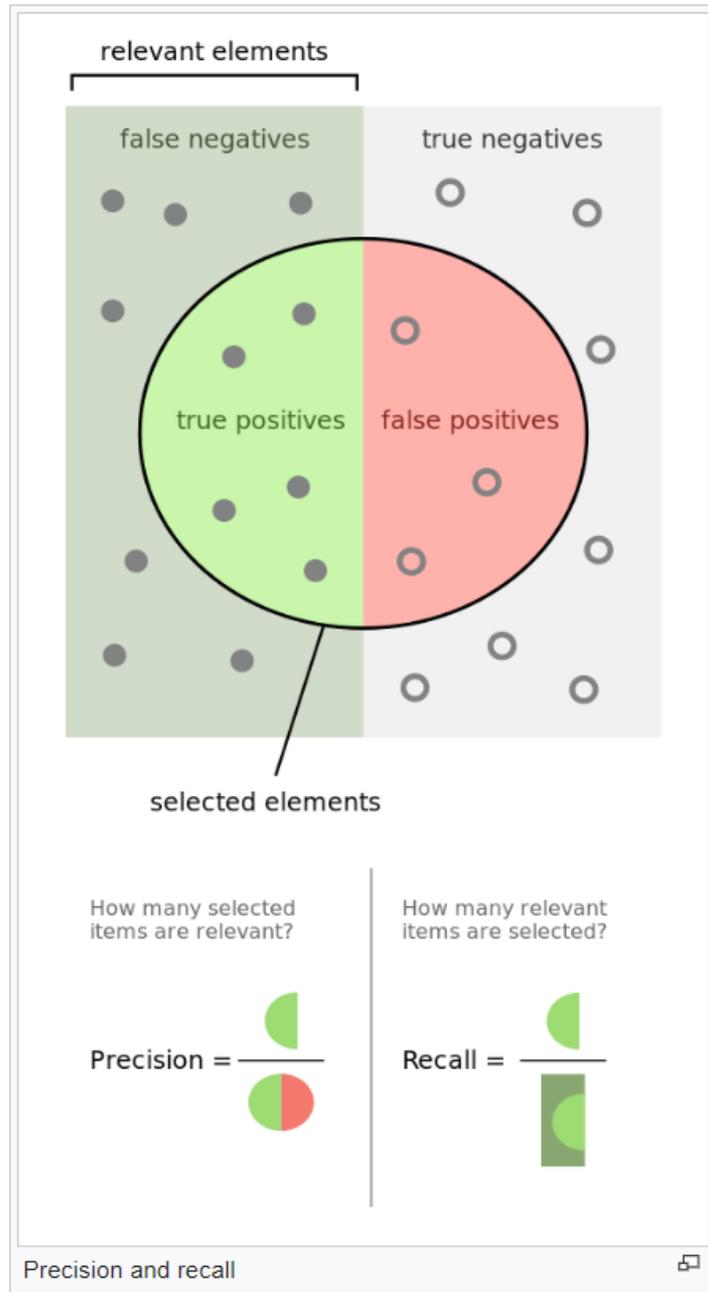
Year 2020 Model Performance (1)

Hourly 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	46	8.61	129.58	0	0.32	0.48	11.38
pure_model_am_day2	32	6.64	70.34	-0.01	0.25	0.72	8.39
pure_model_pm_day2	32	6.48	67.53	-0.01	0.24	0.73	8.22
pure_model_am_day1	34	6.35	65.38	-0.02	0.24	0.74	8.09
pure_model_pm_day1	34	6.31	64.35	-0.02	0.24	0.74	8.02
bias_corrected_am_day1	34	4.27	43.04	0	0.16	0.83	6.56
bias_corrected_pm_day1	34	2.76	27.96	0	0.1	0.89	5.29

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	31	5.68	57.01	0	0.14	0.61	7.55
pure_model_am_day2	20	5.65	49.92	0	0.14	0.66	7.07
pure_model_pm_day2	20	5.26	44.11	0	0.13	0.7	6.64
pure_model_am_day1	24	5.19	42.6	0	0.13	0.71	6.53
pure_model_pm_day1	24	5.05	40.3	0	0.13	0.72	6.35
bias_corrected_am_day1	24	4.78	38.55	0.01	0.12	0.73	6.21
bias_corrected_pm_day1	16	3.09	16.76	0.01	0.08	0.88	4.09



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.

St. Lukes Meridian O3 Site

Year 2020 Model Performance (2)

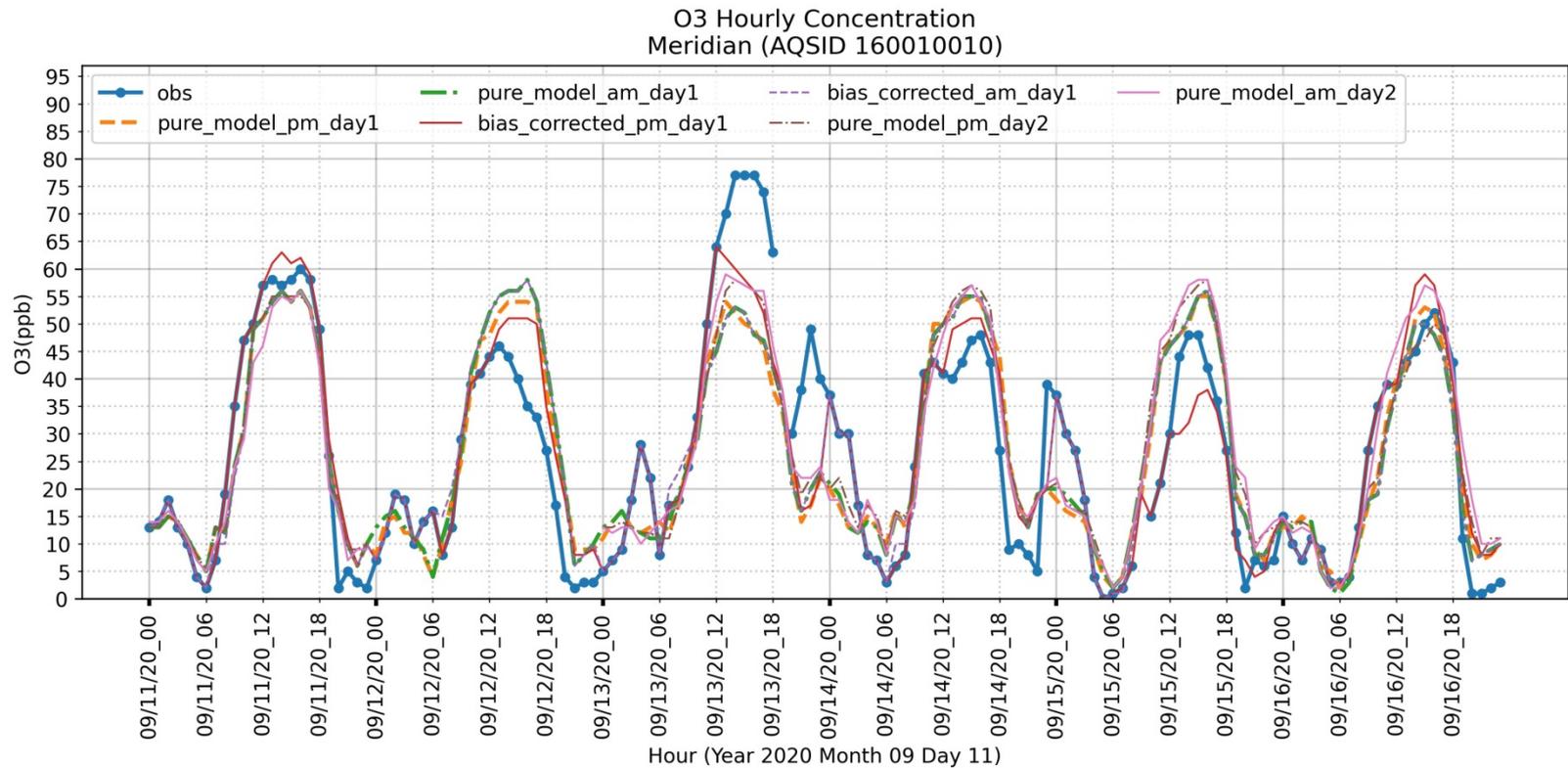
AQI Classification Metrics for All AQI Class

Forecast	accuracy	HSS	KSS
persistence	0.88	0.35	0.36
pure_model_am_day2	0.85	0.3	0.37
pure_model_pm_day2	0.87	0.35	0.41
pure_model_am_day1	0.89	0.38	0.4
pure_model_pm_day1	0.88	0.37	0.39
bias_corrected_am_day1	0.89	0.38	0.4
bias_corrected_pm_day1	0.94	0.68	0.77

AQI Classification Metrics for AQI Class 2 (Yellow)

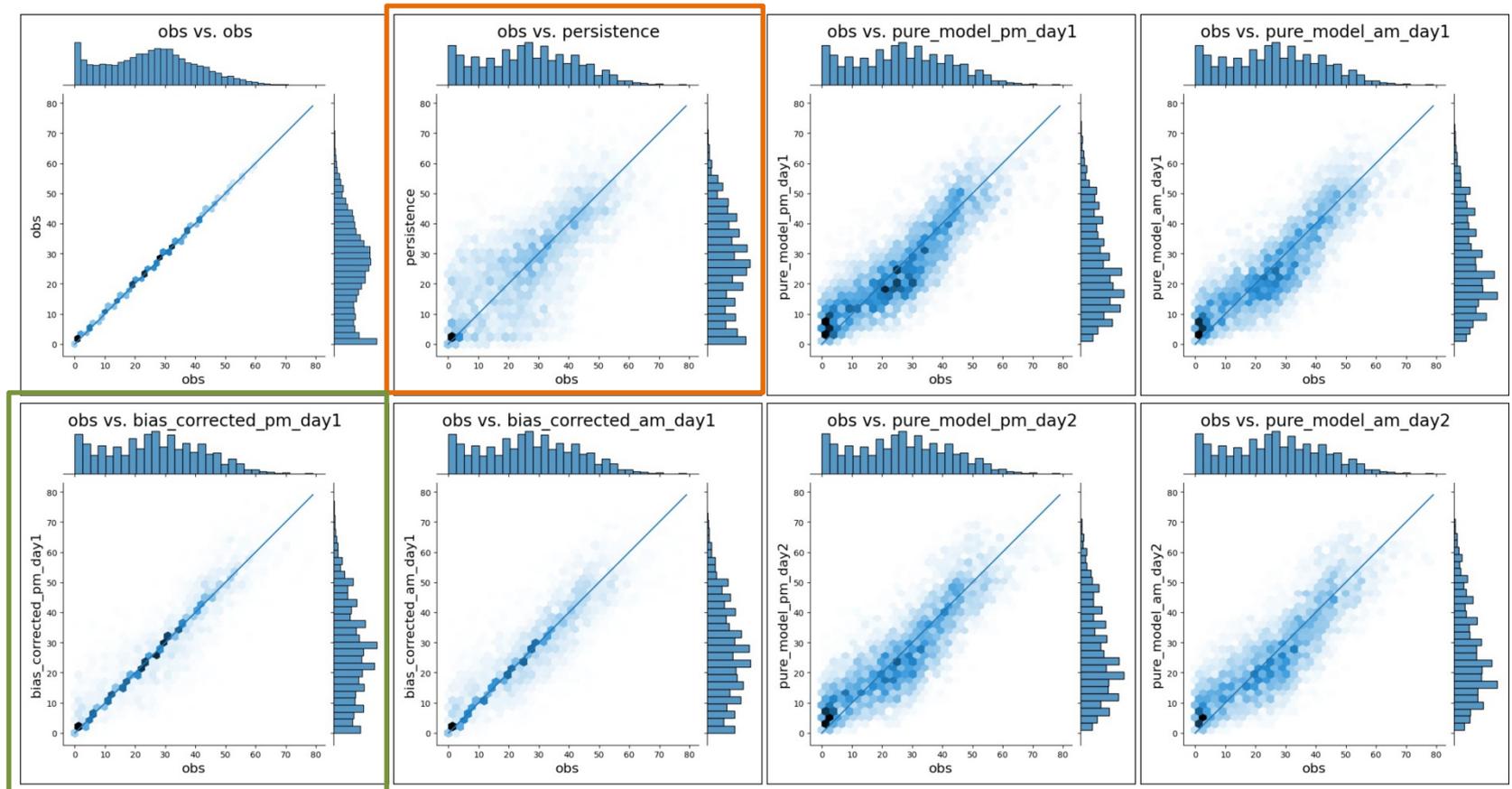
Forecast	AQI_class	precision	recall	f1-score	support
persistence	AQI class 2	0.41	0.44	0.42	25
pure_model_am_day2	AQI class 2	0.32	0.5	0.39	24
pure_model_pm_day2	AQI class 2	0.37	0.52	0.43	25
pure_model_am_day1	AQI class 2	0.43	0.48	0.45	25
pure_model_pm_day1	AQI class 2	0.41	0.48	0.44	25
bias_corrected_am_day1	AQI class 2	0.43	0.48	0.45	25
bias_corrected_pm_day1	AQI class 2	0.62	0.84	0.71	25

St. Lukes Meridian O3 Site Year 2020 Hourly Time Series



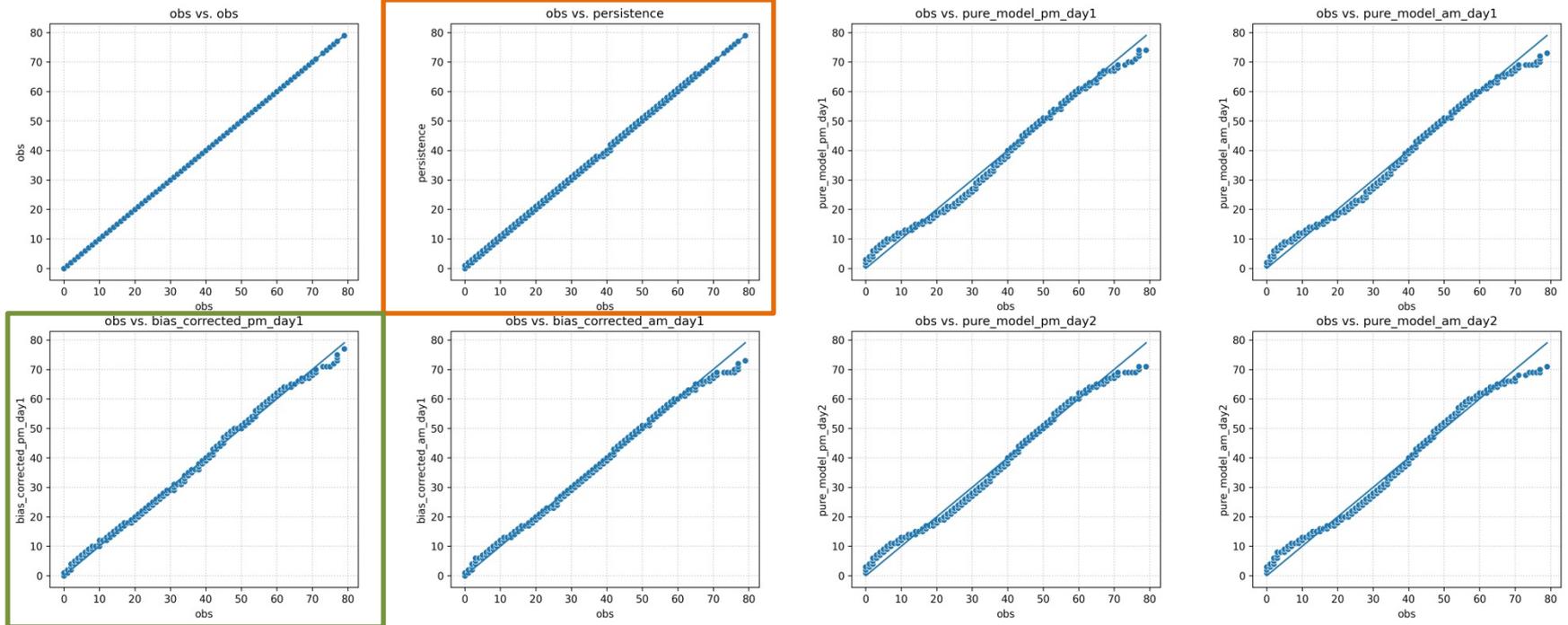
St. Lukes Meridian O3 Site Year 2020 Hourly Joint Plots

O3 Hourly Concentration Joint Plots
Meridian (AQSID 160010010)
Year 2020

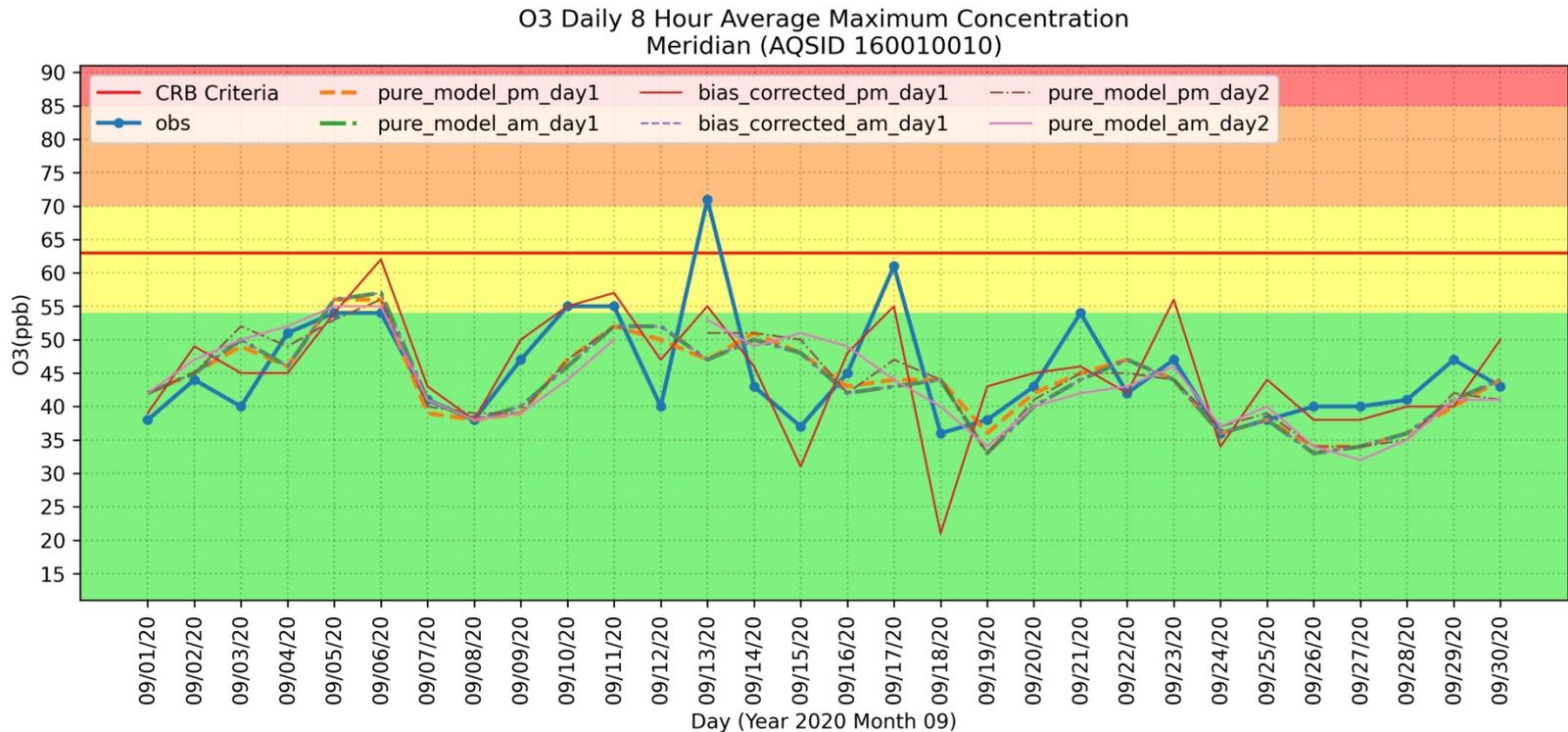


St. Lukes Meridian O3 Site Year 2020 Hourly QQ Plots

O3 Hourly Concentration QQ Plots
Meridian (AQSID 160010010)
Year 2020

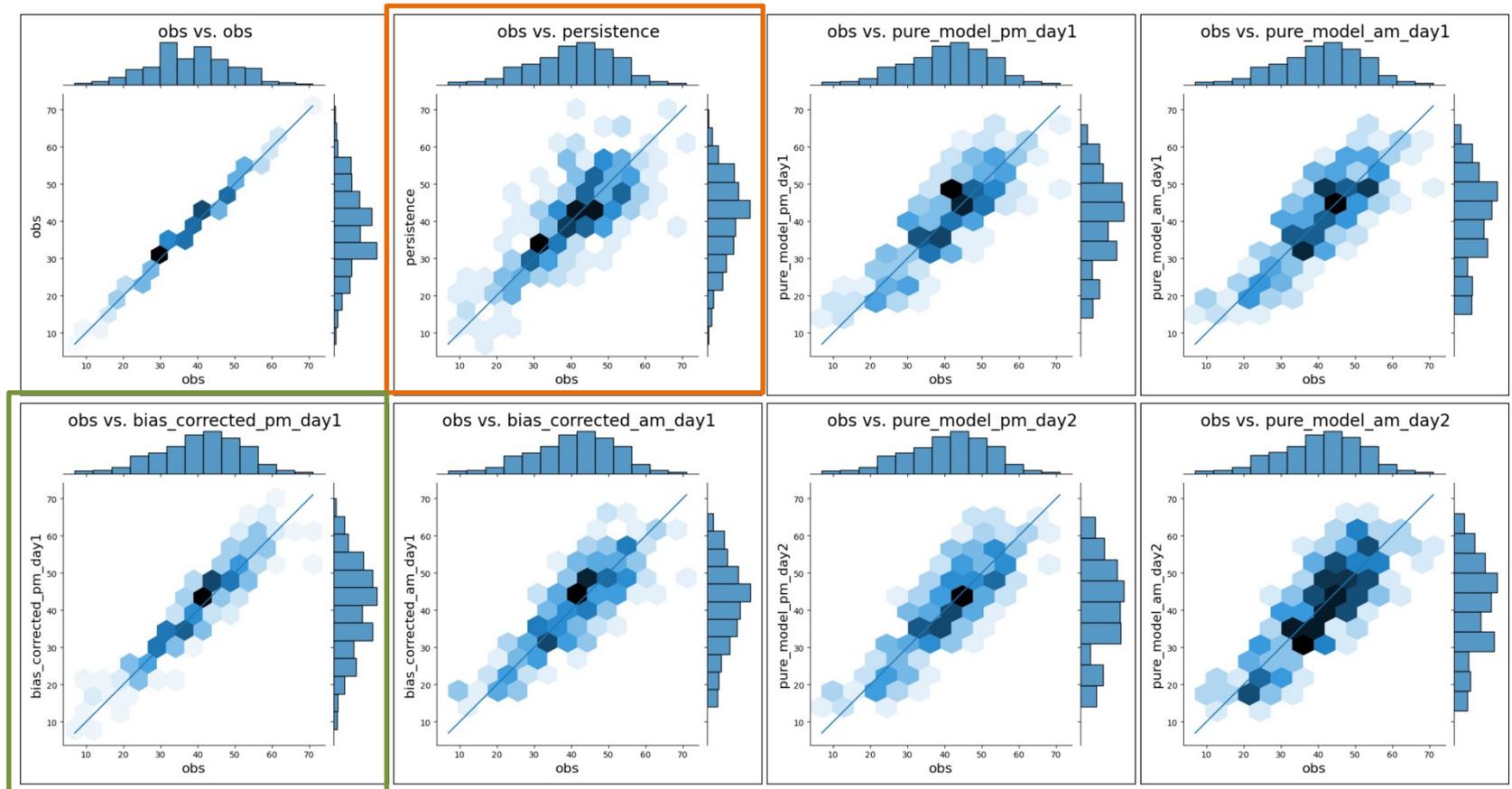


St. Lukes Meridian O3 Site Year 2020 Daily Time Series



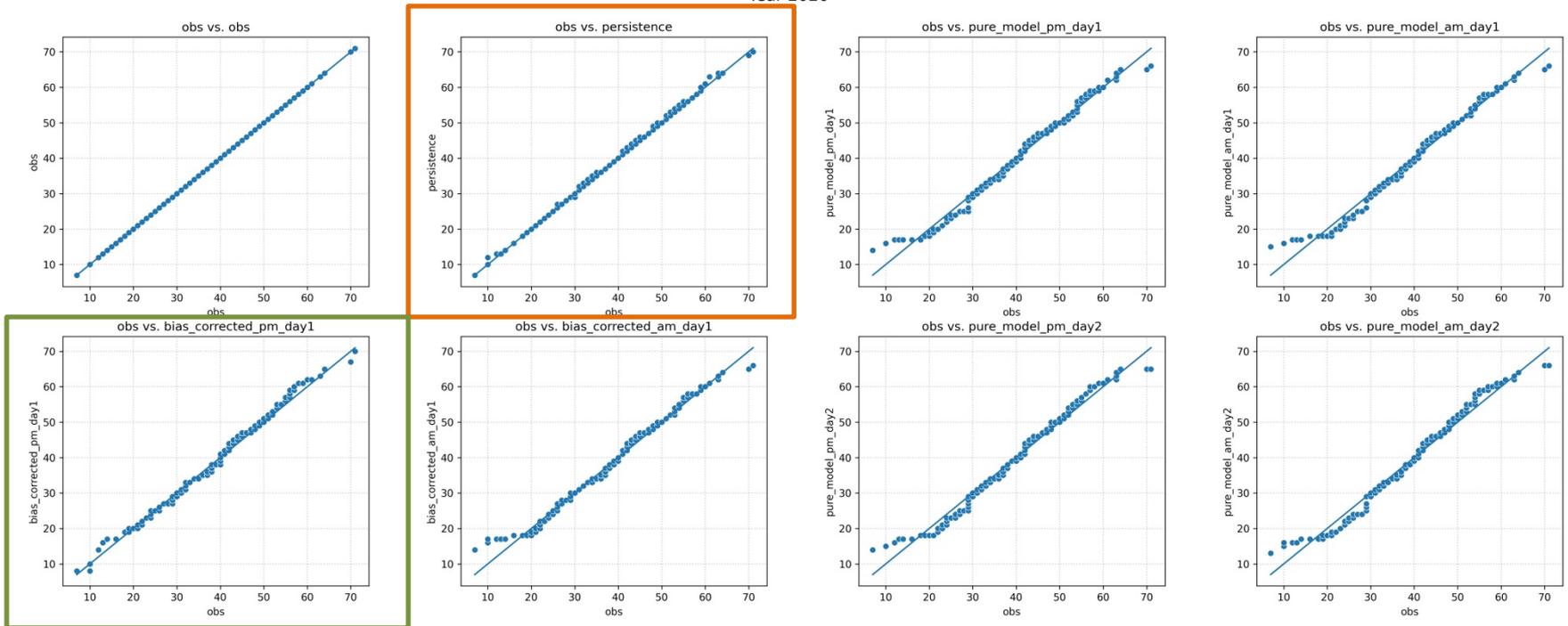
St. Lukes Meridian O3 Site Year 2020 Daily Joint Plots

O3 Daily 8 Hour Average Maximum Concentration Joint Plots
Meridian (AQSID 160010010)
Year 2020



St. Lukes Meridian O3 Site Year 2020 Daily QQ Plots

O3 Daily 8 Hour Average Maximum Concentration QQ Plots
Meridian (AQSID 160010010)
Year 2020

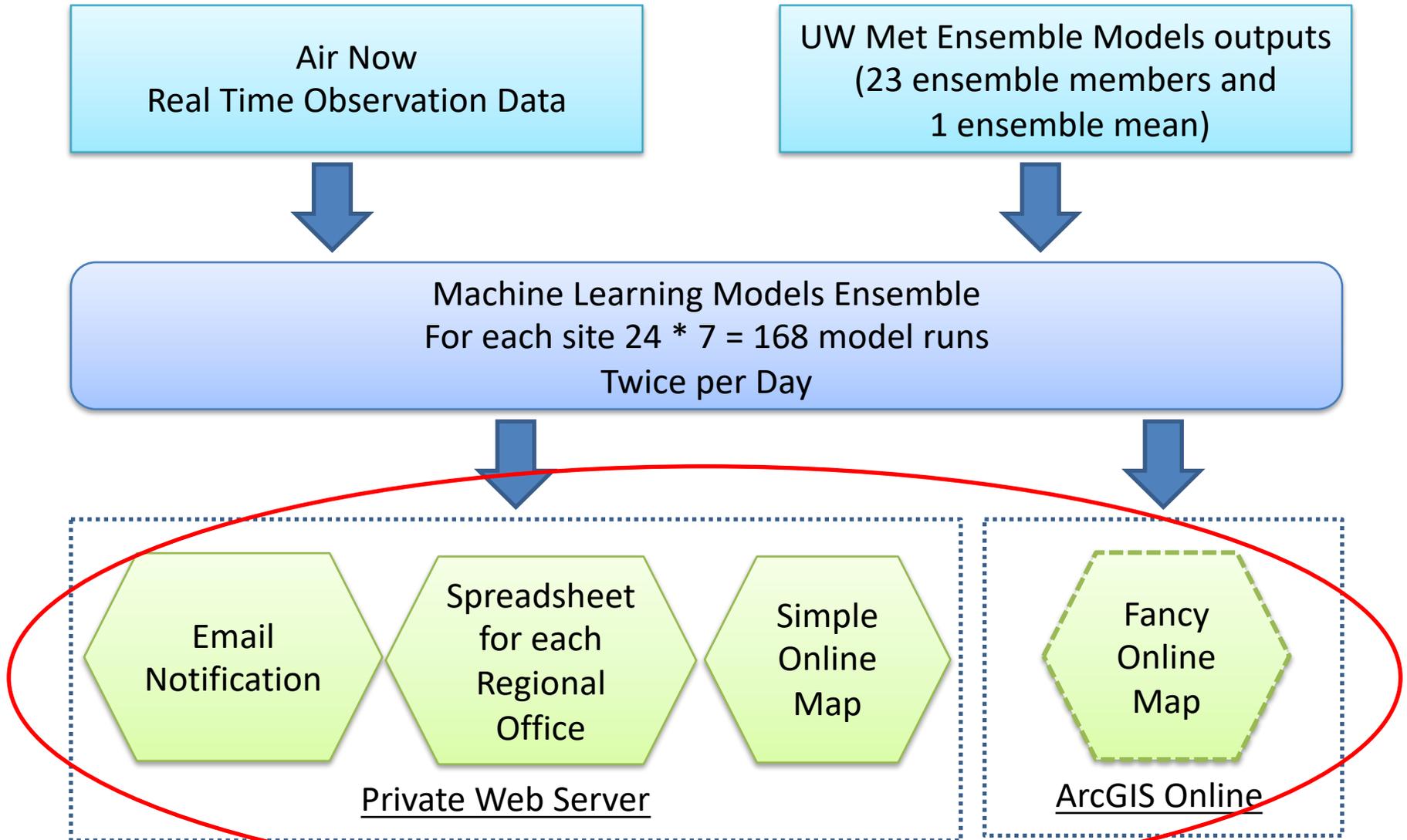


St. Lukes Meridian O3 Site

Year 2020 Model Performance Summary

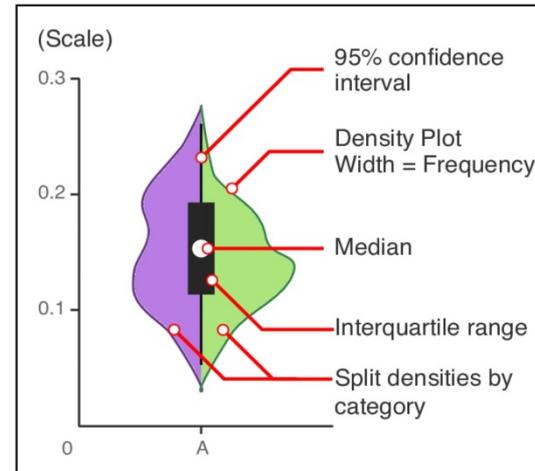
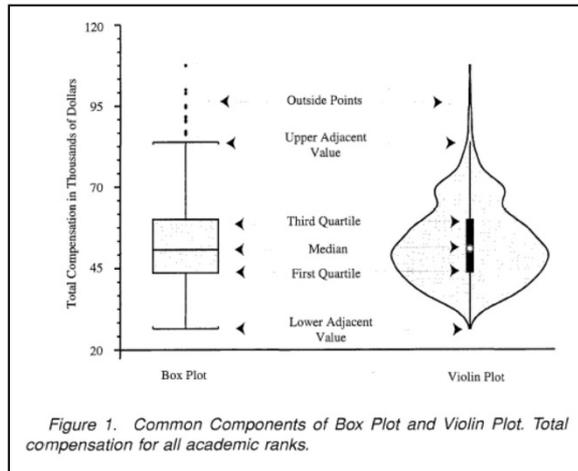
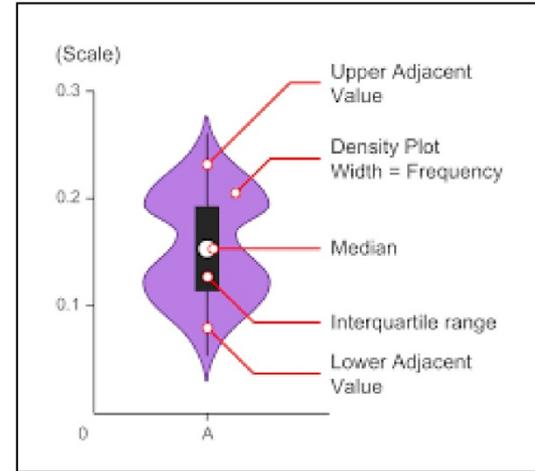
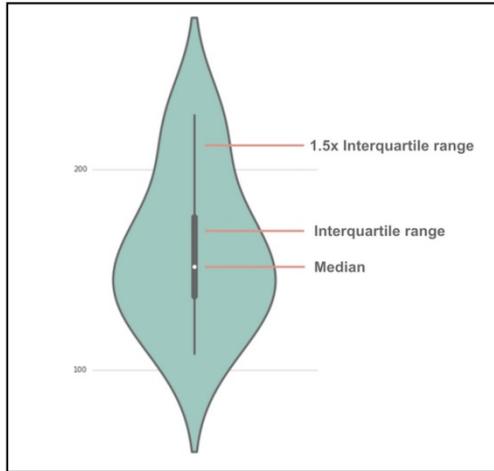
- Better than the persistence method
- Good for Abnormal Detection
- With new data added and model retrained in each year, the model performance will get better and better

Overview of Forecast System



Violin Plot Explained

Present Distribution of Model Ensemble Member Forecasts



Demo

- Email Notification
- Individual Figures
 - Daily 8 Hour Average Maximum Concentration Trend and Detail
 - Daily AQI Trend and Detail
 - Daily Bias (difference) Trend and Detail
 - Hourly Concentration Trend and Detail
- Spreadsheet for each Regional Office
 - Regional offices pick the sites they are interested in
- Simple Online Map

Path Forward

- Expand to PM2.5 sites
- Evaluate the model performance at January of each year
- Retrain the machine learning models with last year's data added at January of each year
- May test it out new features to feed into models in the future

Questions and Discussion

The End

Supplemental Slides

Email Notification

Machine Learning Air Quality Forecast is here! (Forecasted at 12/21/2020 Afternoon)

Wei.Zhang@deq.idaho.gov

Sent: Mon 12/21/2020 1:27 PM

To: Wei Zhang

Wei,

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

http://10.220.98.54/ml_forecast_outputs/2020/20201221_pm/

OR through simple online map http://10.220.98.54/ml_forecast_outputs/2020/20201221_pm/20201221_pm_forecast_0Map.html

To access forecasts for previous days, please go to:

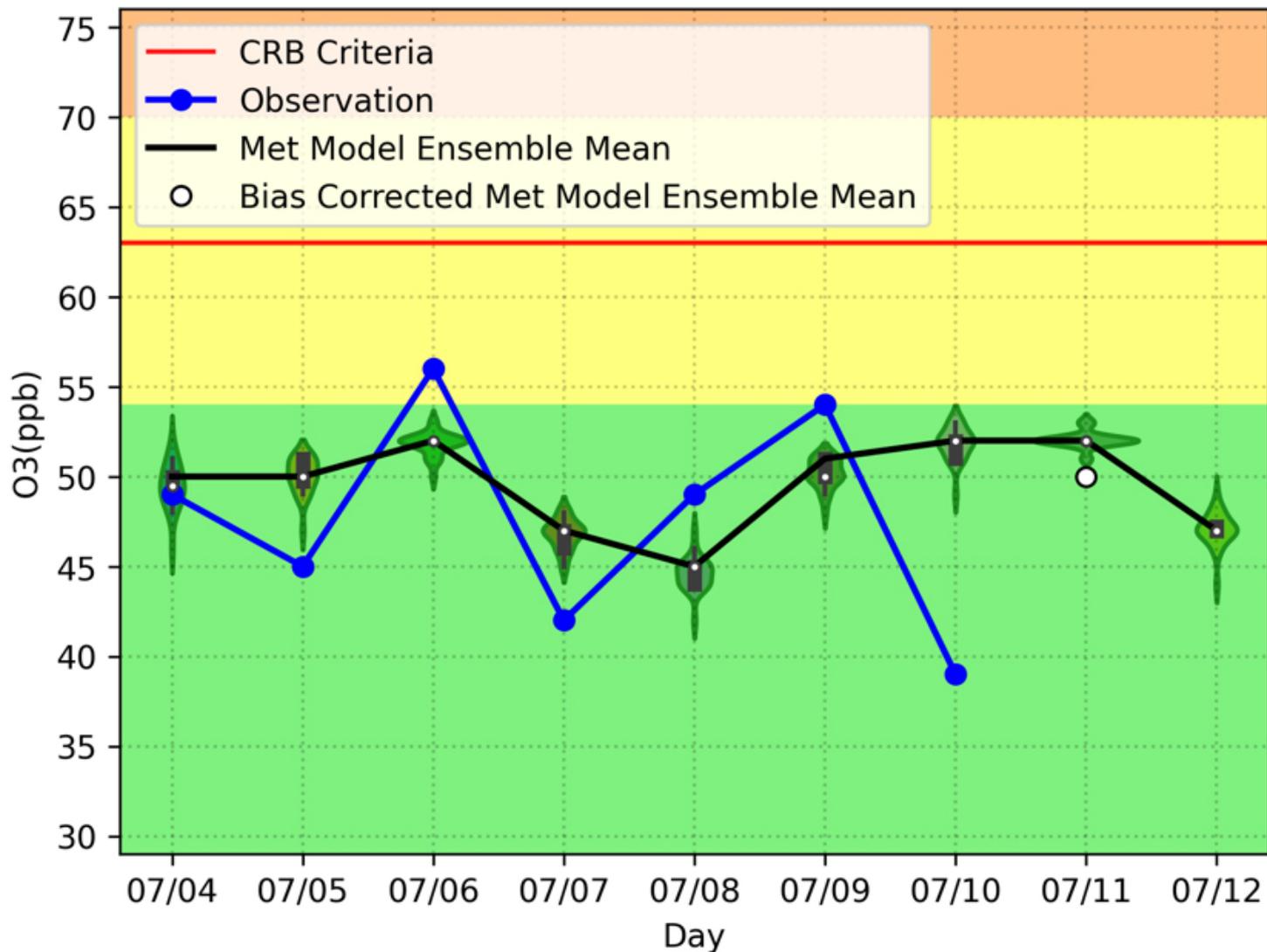
http://10.220.98.54/ml_forecast_outputs/

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.

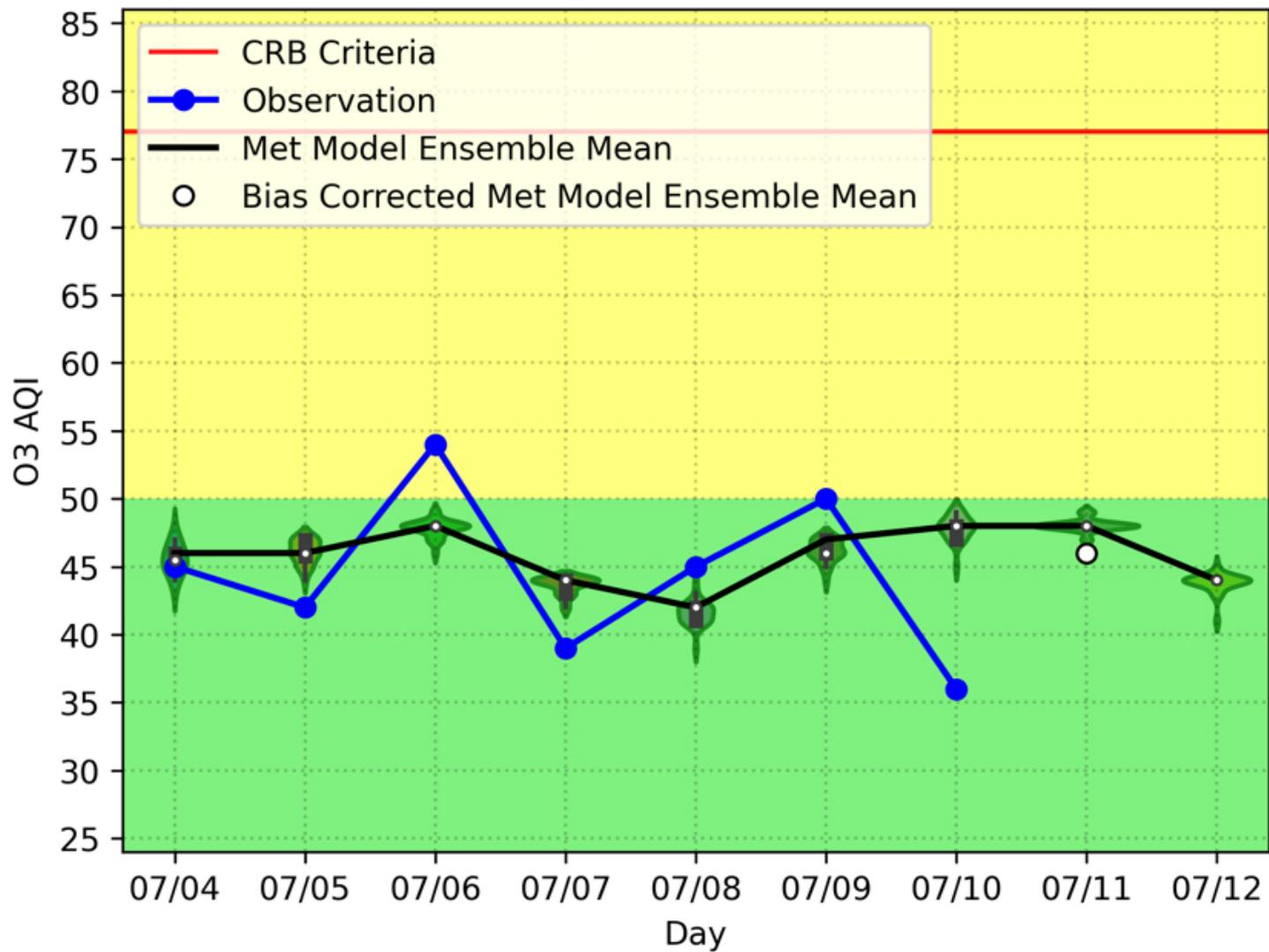
Hopefully this tool makes your life easier. Please provide us any feedback you may have.

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

O3 Daily 8 Hour Average Maximum Concentration
Meridian (AQSID 160010010)
Forecasted at 07/11/2020 Afternoon

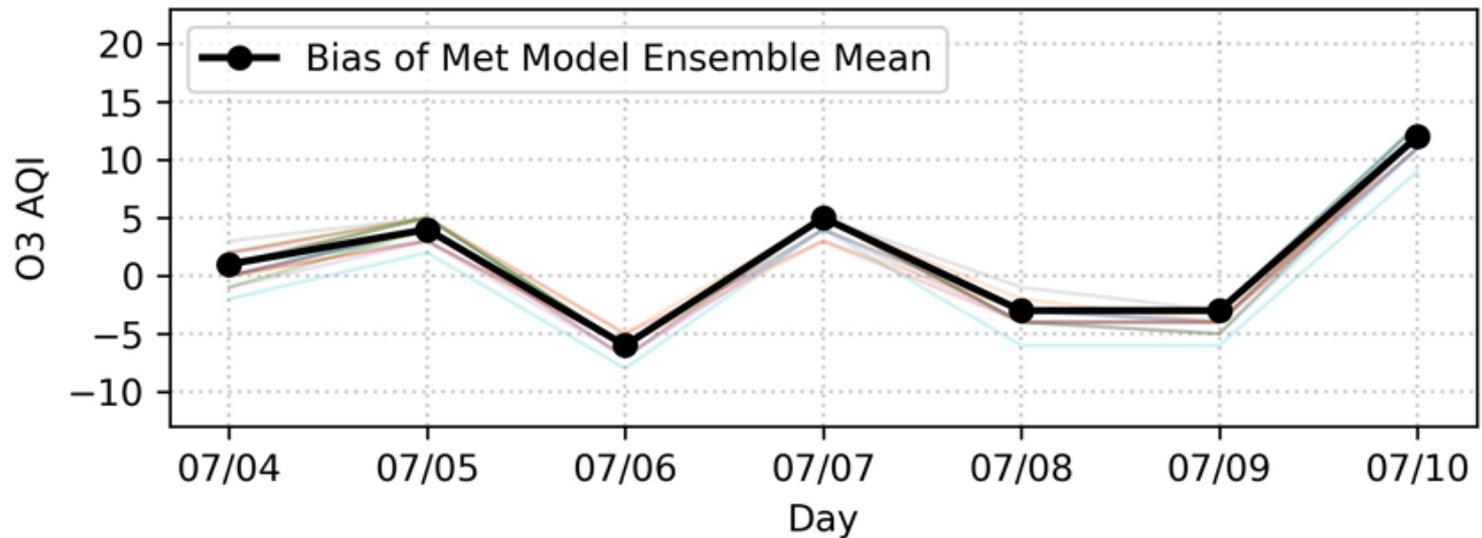
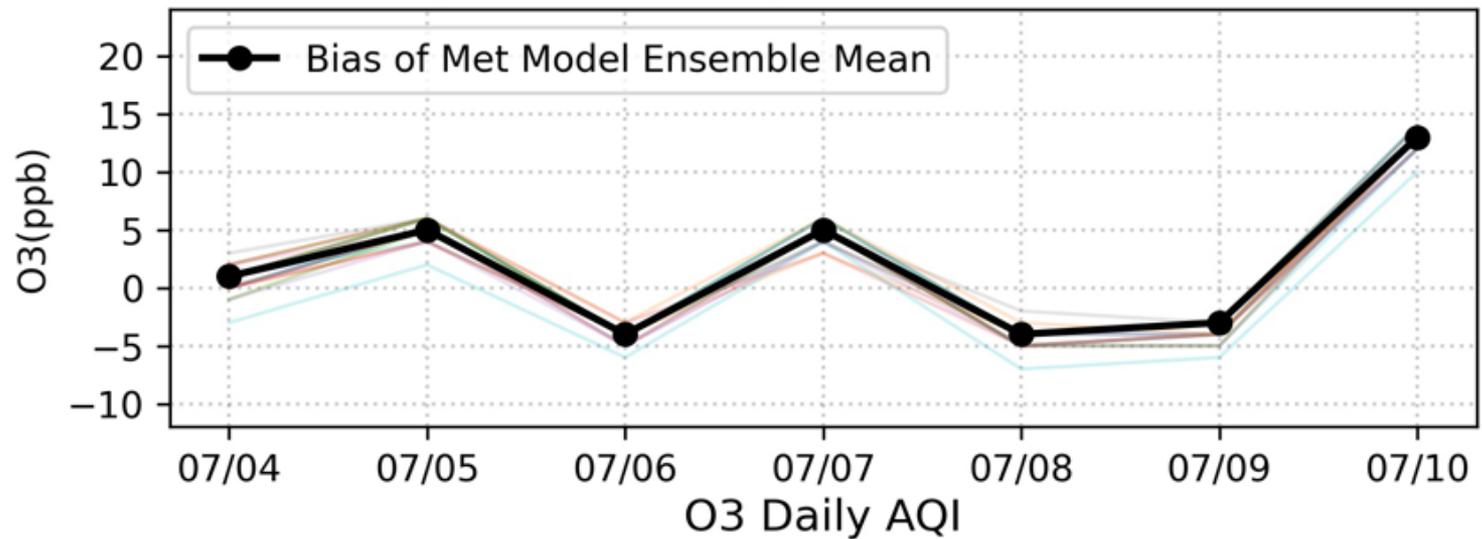


O3 Daily AQI
Meridian (AQSID 160010010)
Forecasted at 07/11/2020 Afternoon

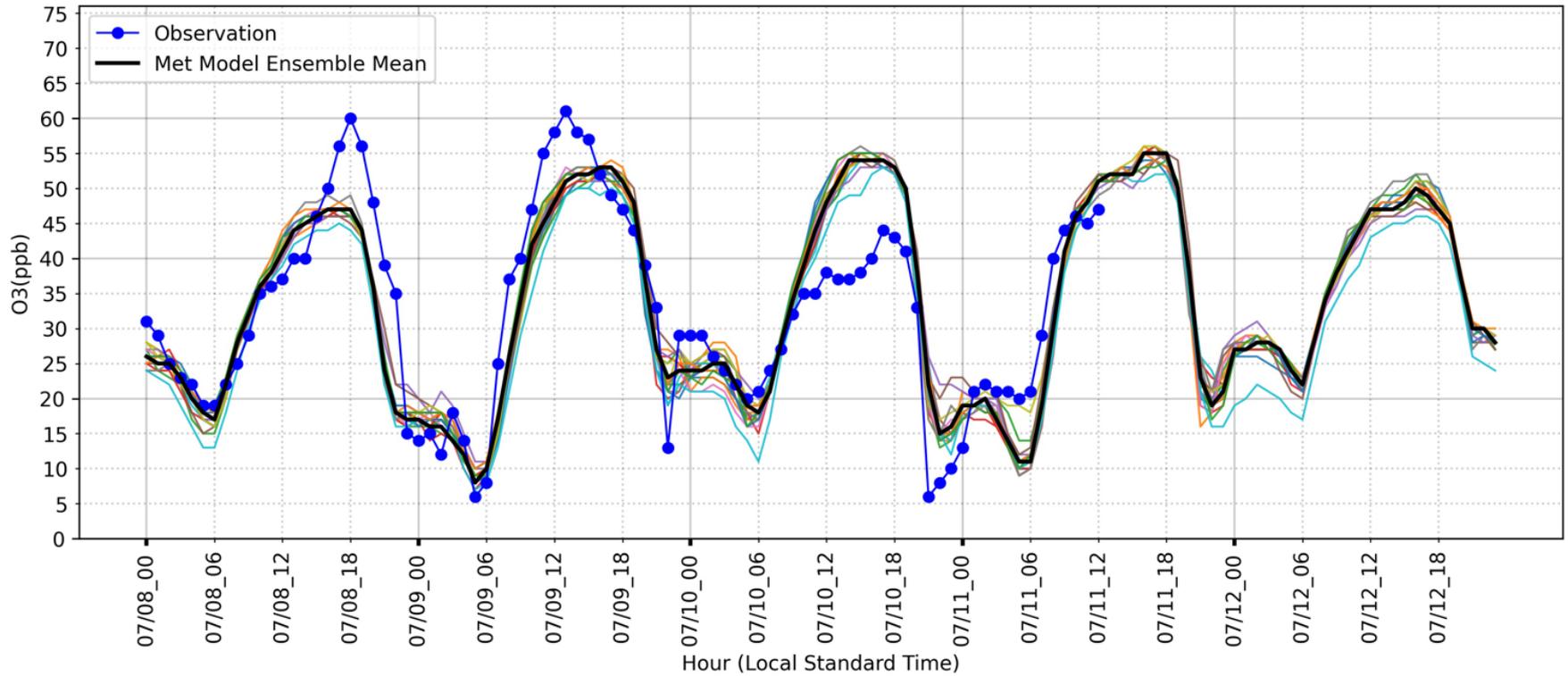


O3 Daily Forecast Bias Meridian (AQSID 160010010)

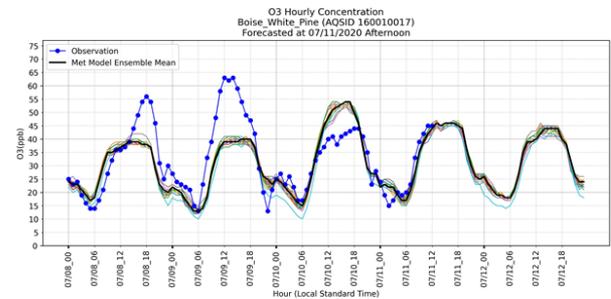
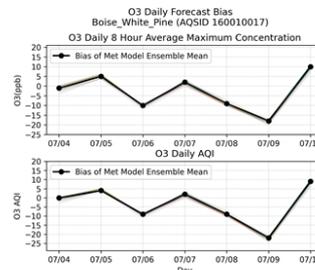
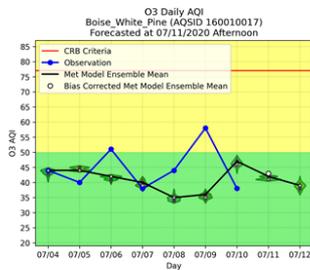
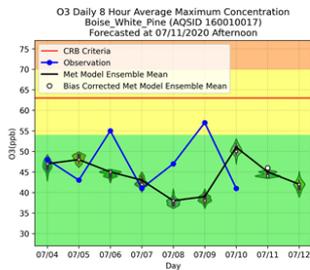
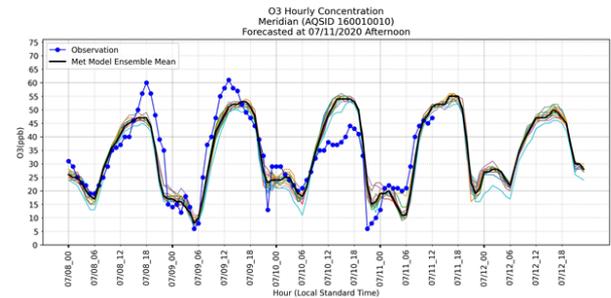
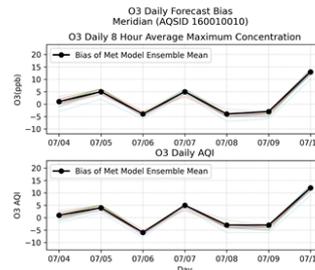
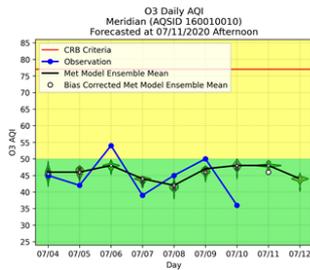
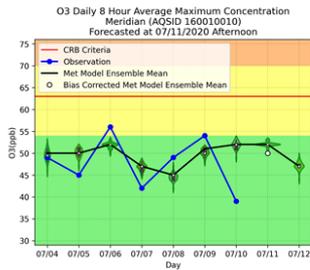
O3 Daily 8 Hour Average Maximum Concentration



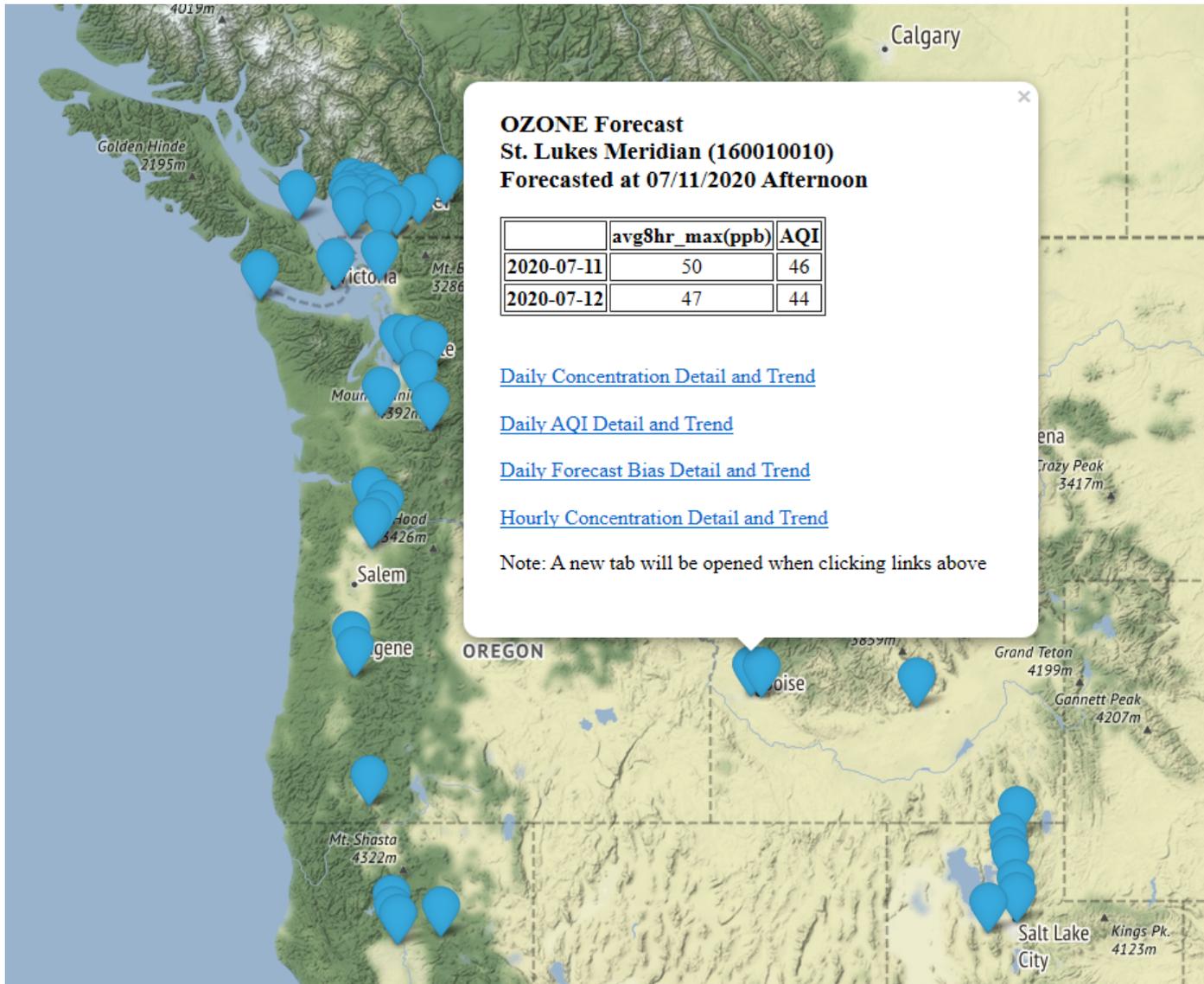
O3 Hourly Concentration
Meridian (AQSID 160010010)
Forecasted at 07/11/2020 Afternoon



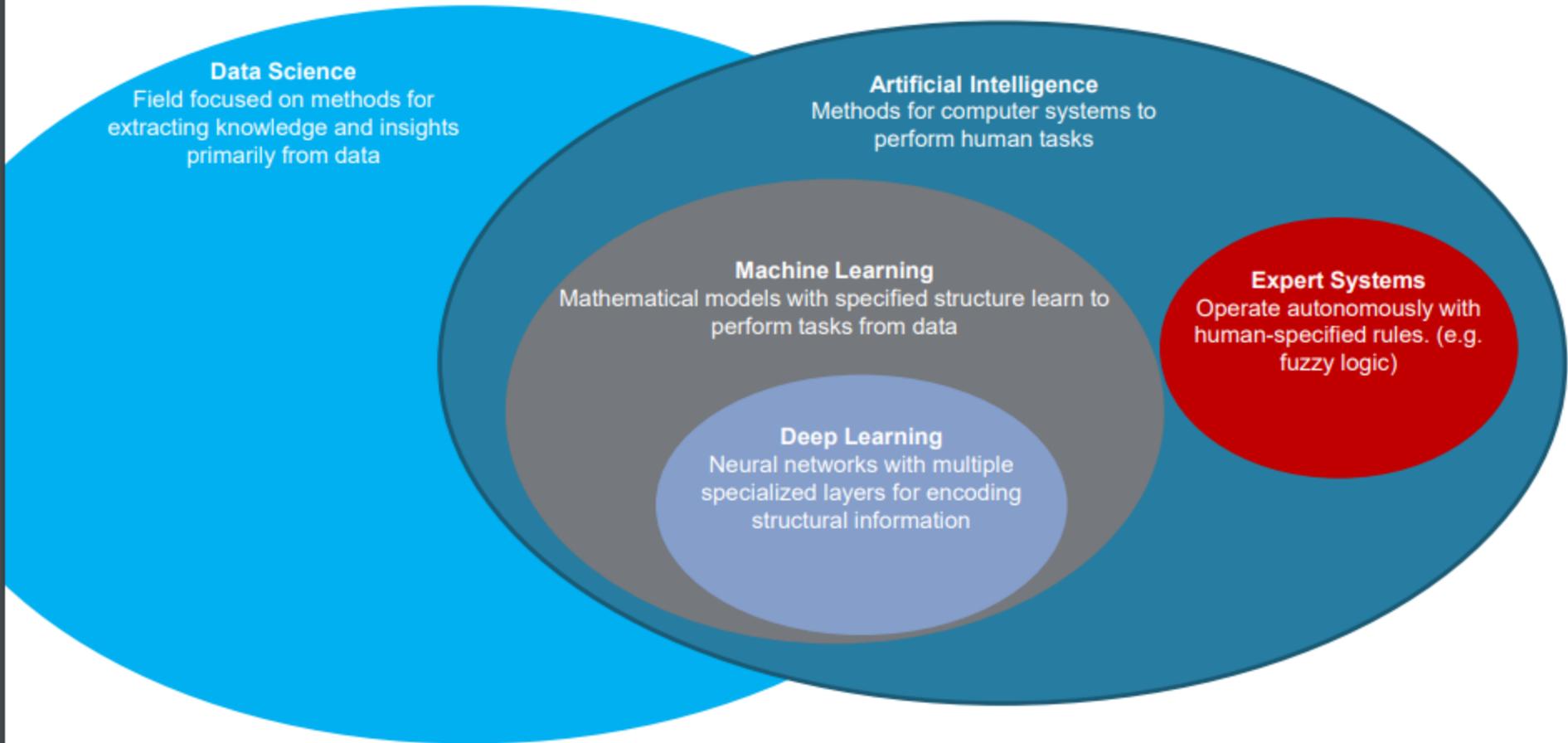
Spreadsheet per Regional Office



Simple Online Map

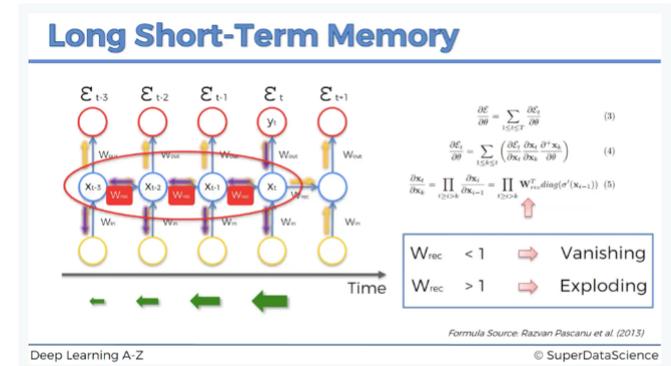
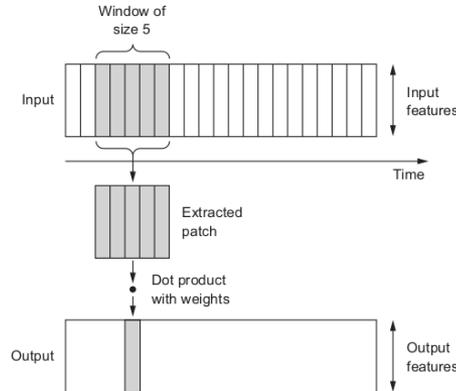
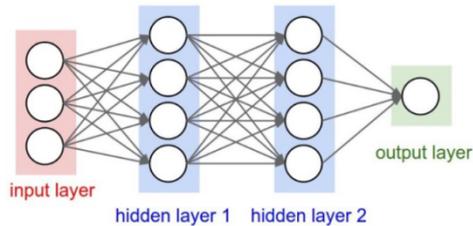


The Data Science Taxonomy



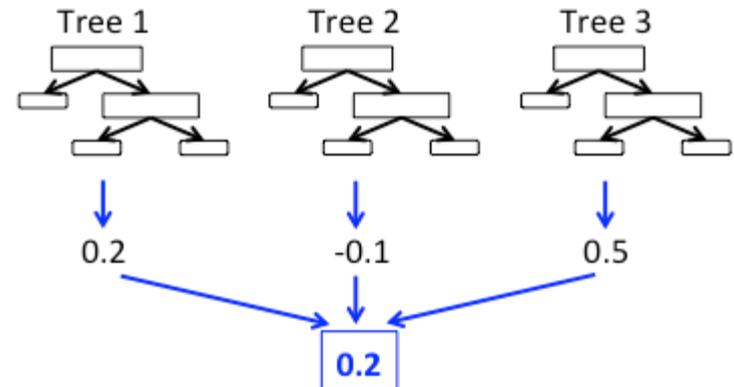
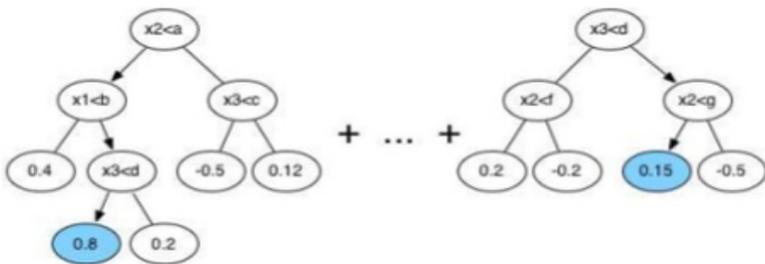
Neural Network

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



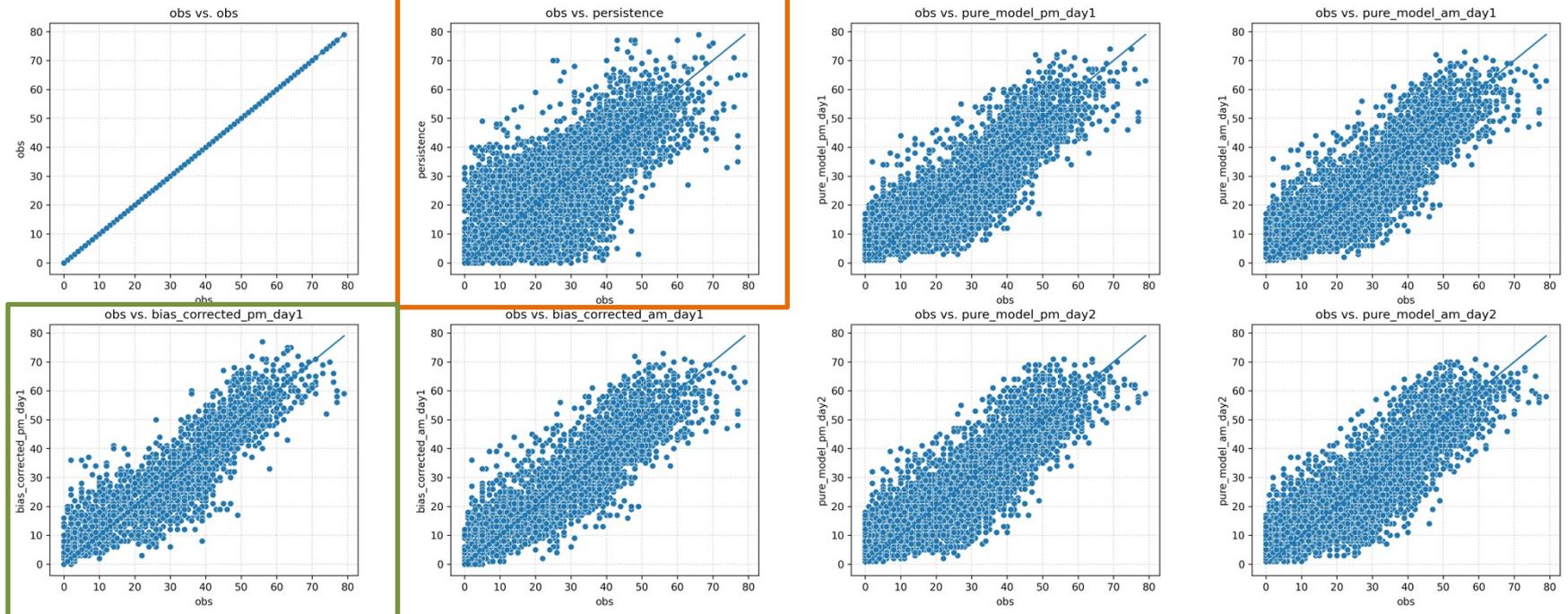
Tree based Methods (XGBoost)

- Pure 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



St. Lukes Meridian O3 Site Year 2020 Hourly Scatter Plots

O3 Hourly Concentration Scatter Plots
Meridian (AQSID 160010010)
Year 2020



St. Lukes Meridian O3 Site Year 2020 Daily Scatter Plots

O3 Daily 8 Hour Average Maximum Concentration Scatter Plots
Meridian (AQSID 160010010)
Year 2020

