

A Machine Learning Method for Estimating Atmospheric Trace Gas Concentration Baselines

Kirstin Gerrand^{1, 2}, Elena Fillola^{1, 3}, Alistair J. Manning^{4, 1}, Jgor Arduini⁵, Paul B. Krummel⁶, Chris R. Lunder⁷, Jens Mühle⁸, Simon O'Doherty¹, Sunyoung Park⁹, Ronald G. Prinn¹⁰, Stefan Reimann¹¹, Dickon Young¹, and Matthew Rigby¹

¹Atmospheric Chemistry Research Group, School of Chemistry, University of Bristol, Bristol, UK

²New Zealand Institute for Earth Science, Wellington, New Zealand

³School of Engineering Mathematics, University of Bristol, Bristol, UK

⁴Hadley Centre, Met Office, Exeter, UK

⁵Department of Pure and Applied Sciences, University of Urbino, Urbino, 61029, Italy

⁶CSIRO Environment, Aspendale, VIC, Australia

⁷NILU, Kjeller, Norway

⁸Scripps Institution of Oceanography, University of California San Diego, La Jolla, CA, USA

⁹Kyungpook Institute of Oceanography, Kyungpook National University, Daegu, Republic of Korea

¹⁰Center for Sustainability Science and Strategy, Massachusetts Institute of Technology, Cambridge, MA, USA

¹¹Laboratory for Air Pollution/Environmental Technology, Empa, Dübendorf, Switzerland

Supplementary Material

Here, we show additional figures and tables that elaborate on model choices and performance.

5 An example comparison of ECMWF and UK Met Office wind speed and direction at Mace Head, Ireland is shown in Figure S1, indicating that differences between the meteorology used to drive the NAME model (UK Met Office analysis) and that used in training our algorithm (ECMWF ERA5) are small and unlikely to contribute substantially to algorithm performance issues.

A summary of the meteorological parameters used as features in the machine learning model are provided in Table S1.

Metrics derived from the Random Forest (RF) model are provided in Table S2, based on hyperparameter sets in Tables S3a and S3b, and the the RF confusion matrix and heatmap is shown in Figs. S2 and S3, respectively.

10 Hyperparameter sets for the final MLP models are shown in Tables S4a and S4b. The feature importance for the MLP models are shown in Table S5. Plots indicating MLP model performance for a range of species and sites are shown in Section 4.3.

1 Intercomparison of ECMWF and Met Office UM Meteorological Fields

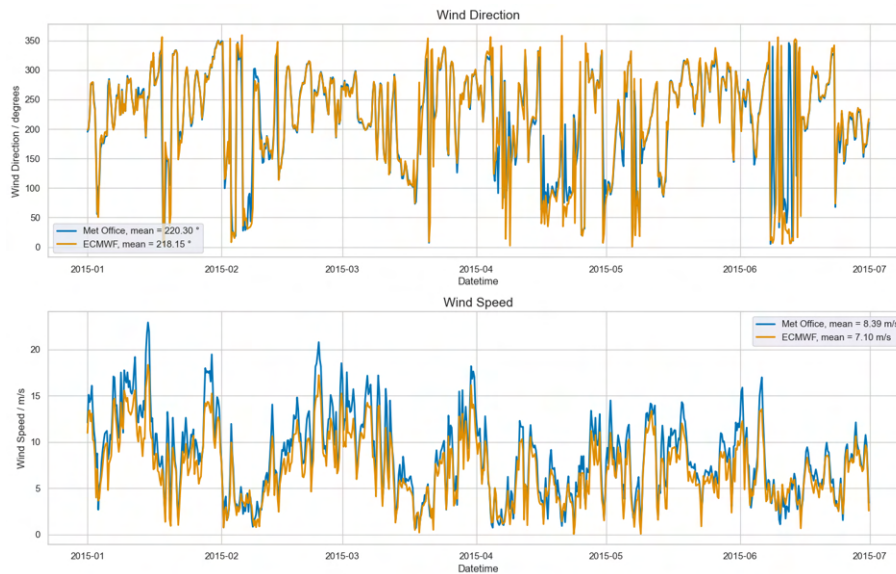


Figure S1. Subplots comparing ECMWF and Met Office UM 10m winds across a six-month sample period (Jan-June 2015) at Mace Head, Ireland. Top figure compares wind directions, and bottom, wind speeds. Legend contains average values for additional quantitative comparison.

2 Model Inputs

Table S1. Input variable categories used in the ML models. Each variable is extracted from 17 horizontally distributed points around the site; the data collection site itself and 16 from surrounding nearest grid points based on two 3x3 grids spanning $\pm 5^\circ$ and $\pm 10^\circ$ latitude and longitude.

<i>Category</i>	<i>Units</i>
Boundary layer height	m
Surface Pressure	Pa
10m u-wind	m s^{-1}
10m u-wind (T-6)	m s^{-1}
500hPa u-wind	m s^{-1}
500hPa u-wind (T-6)	m s^{-1}
850hPa u-wind	m s^{-1}
850hPa u-wind (T-6)	m s^{-1}
10m v-wind	m s^{-1}
10m v-wind (T-6)	m s^{-1}
500hPa v-wind	m s^{-1}
500hPa v-wind (T-6)	m s^{-1}
850hPa v-wind	m s^{-1}
850hPa v-wind (T-6)	m s^{-1}
Time of Day	
Day of Year	

3 Random Forest Summary

15 3.1 Final Random Forest Model Outcomes

Table S2. A tabular summary of the final RF model outcomes, showing precision, recall and F1 score values for each of the AGAGE sites. Scores can range from 0 to 1, with higher scores indicating higher performance.

<i>Station Designation</i>	<i>Precision</i>	<i>Recall</i>	<i>F1 score</i>
CGO	0.953	0.818	0.880
CMN	0.957	0.588	0.728
GSN	0.925	0.608	0.734
JFJ	0.941	0.674	0.785
MHD	0.917	0.672	0.776
RPB	0.884	0.521	0.655
SMO	0.896	0.404	0.557
THD	0.915	0.437	0.591
ZEP	0.948	0.718	0.817

3.2 Final Hyperparameter Sets

Note that hyperparameters not listed are as default in the documentation.

Table S3. Hyperparameter sets used in the final versions of the RF models for AGAGE sites.

(a) Four AGAGE sites: Kennaoak/Cape Grim (CGO), Monte Cimone (CMN), Gosan (GSN), Jungfraujoch (JFJ).

<i>Parameter</i>	<i>CGO</i>	<i>CMN</i>	<i>GSN</i>	<i>JFJ</i>
random_state	42	42	42	42
n_estimators	100	100	200	100
max_depth	5	3	5	5
criterion	'entropy'	'entropy'	'entropy'	'entropy'
bootstrap	False	False	False	False
max_features	'sqrt'	'sqrt'	None	'sqrt'

(b) Five AGAGE sites: Mace Head (MHD), Ragged Point (RPB), Cape Matatula (SMO), Trinidad Head (THD), Zeppelin (ZEP).

<i>Parameter</i>	<i>MHD</i>	<i>RPB</i>	<i>SMO</i>	<i>THD</i>	<i>ZEP</i>
random_state	42	42	42	42	42
n_estimators	200	100	50	100	200
max_depth	8	5	5	5	5
criterion	'entropy'	'entropy'	'log_loss'	'gini'	'entropy'
bootstrap	False	False	False	True	False
min_samples_split	5	5	3	2	2
max_features	'sqrt'	'log2'	'sqrt'	'sqrt'	'sqrt'

3.3 Random Forest Confusion Matrix Map

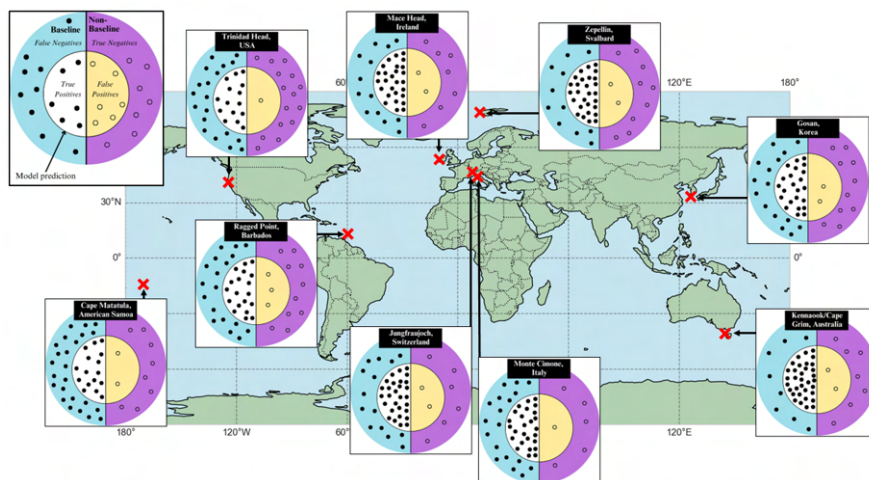


Figure S2. A map showing the locations of the nine AGAGE sites, with confusion matrix-derived plots at each location. Each confusion matrix was normalised, to remove visual differences due to testing set sizes; each point represents approximately 2% of the total test set (rounding means that the total number of points on each plot range from 49 to 51). The left half of each circle represents true baseline points, and the right half true non-baseline points. The inner circle shows the random forest model prediction of baseline points. The key in the top left indicates where true or false positives and negatives lie, as explained in the main text.

3.4 Random Forest MAPE Heatmap

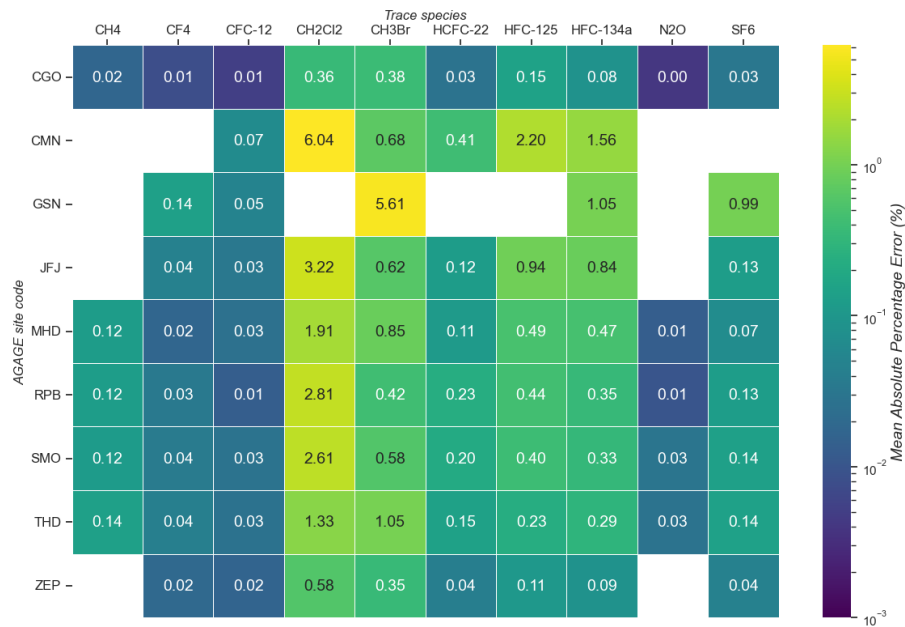


Figure S3. Mean Absolute Percentage Errors (MAPE) for the random forest model for selected AGAGE species across the nine monitoring sites. White denotes species that were not available at a particular site.

4.1 Final Hyperparameter Sets

Note that hyperparameters not listed are as default in the documentation.

Table S4. Hyperparameter sets used in the final versions of the MLP models for AGAGE sites.

(a) Four AGAGE sites: Kennaoook/Cape Grim (CGO), Monte Cimone (CMN), Gosan (GSN), Jungfraujoch (JFJ).

<i>Parameter</i>	<i>CGO</i>	<i>CMN</i>	<i>GSN</i>	<i>JFJ</i>
random_state	42	42	42	42
max_iter	1000	1000	1000	50
hidden_layer_sizes	(50,)	(100,)	(100,)	(25,)
shuffle	False	False	False	False
alpha	0.05	0.5	0.0001	0.1
learning_rate	'constant'	'constant'	'constant'	'invscaling'
batch_size	100	250	100	100
early_stopping	True	False	False	True
learning_rate_init	0.001	0.01	0.01	0.01
epsilon	1e-8	1e-8	1e-8	2e-10
beta_2	0.9	0.9	0.9	0.8
tol	1e-4	1e-4	1e-4	1e-3

(b) Five AGAGE sites: Mace Head (MHD), Ragged Point (RPB), Cape Matatula (SMO), Trinidad Head (THD), Zeppelin (ZEP).

<i>Parameter</i>	<i>MHD</i>	<i>RPB</i>	<i>SMO</i>	<i>THD</i>	<i>ZEP</i>
random_state	42	1	42	42	42
max_iter	1000	500	1000	1000	1000
hidden_layer_sizes	(200,)	(200,)	(100,)	(100,)	(200,150,)
shuffle	False	False	False	False	False
alpha	0.05	0.0001	0.001	0.05	0.05
learning_rate	'constant'	'constant'	'constant'	'constant'	'constant'
batch_size	100	100	100	100	100
early_stopping	True	True	True	True	True
learning_rate_init	0.01	0.001	0.001	0.001	0.01
beta_2	0.9	0.9	0.9	0.9	0.9

4.2 Feature Importance

Table S5. A tabular summary of the top three most influential feature groups when training the final MLP model for each site. u- and v- wind encompass all heights, locations and times. 'PBLH' represents planetary boundary layer height.

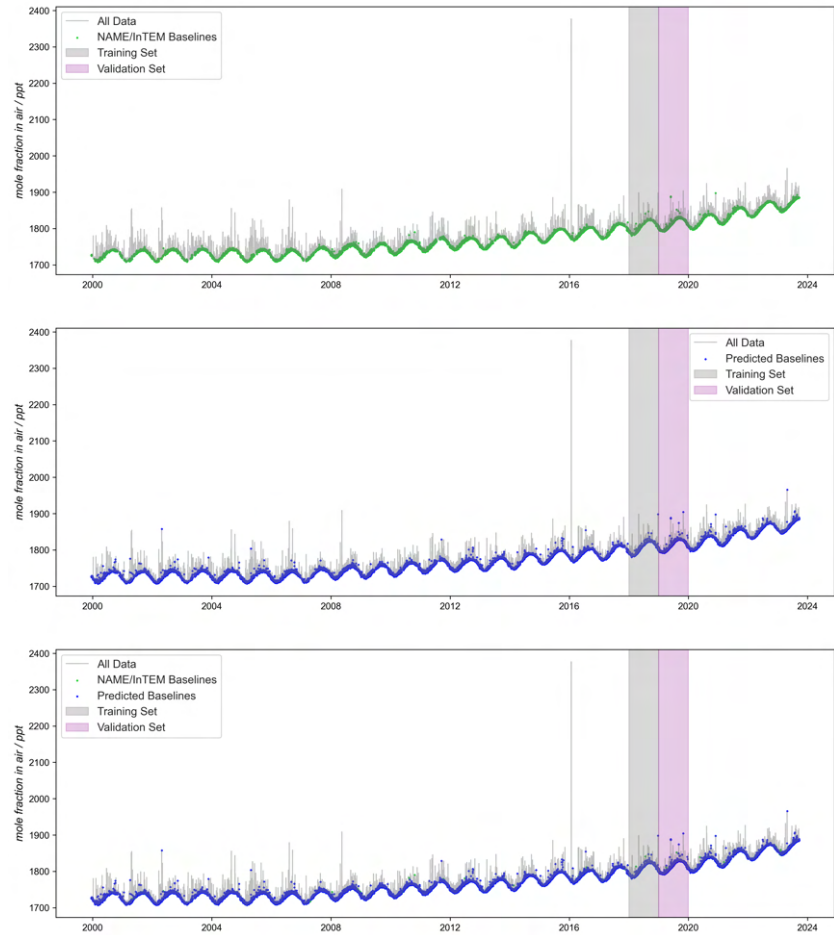
<i>Station Designation</i>	<i>Most Importance</i>	<i>Second</i>	<i>Third</i>
CGO	u-wind	v-wind	PBLH
CMN	PBLH	u-wind	Surface pressure
GSN	PBLH	u-wind	Surface pressure
JFJ	v-wind	u-wind	Surface Pressure
MHD	u-wind	v-wind	PBLH
RPB	u-wind	v-wind	PBLH
SMO	u-wind	w-wind	PBLH
THD	PBLH	v-wind	u-wind
ZEP	v-wind	u-wind	PBLH

4.3 MLP Plots

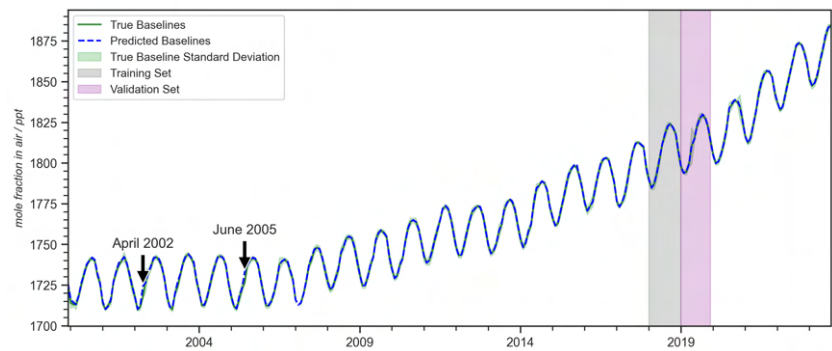
25 4.4 Kennaoook/Cape Grim, Australia

4.4.1 CH₄

Mole fraction time series

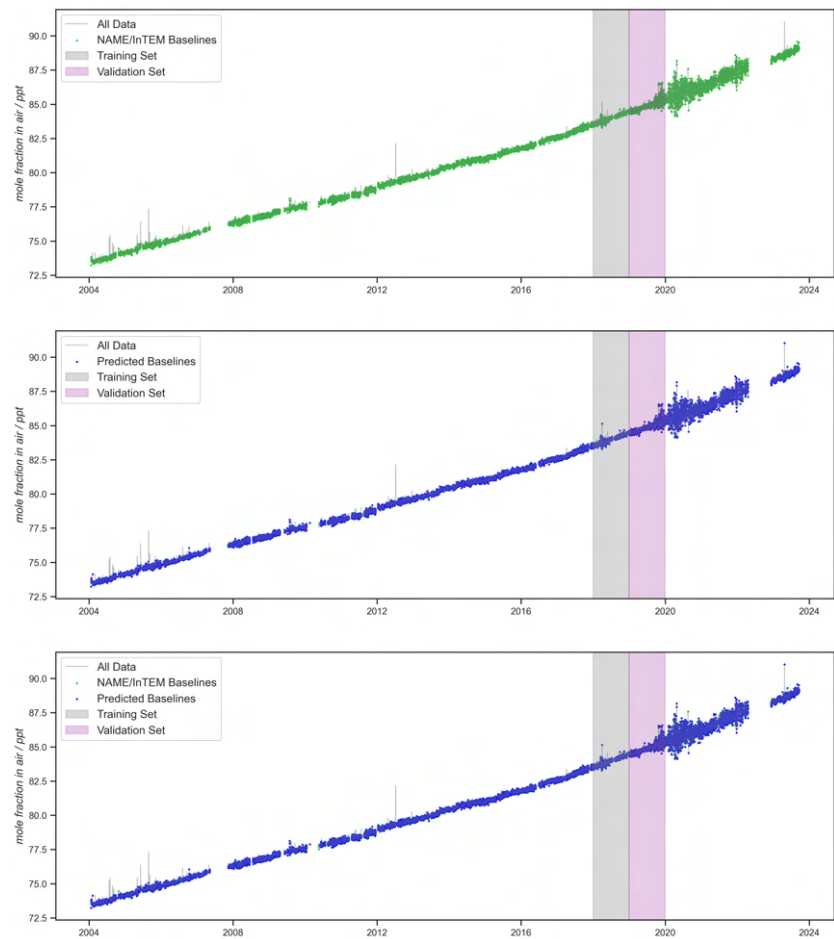


Monthly means

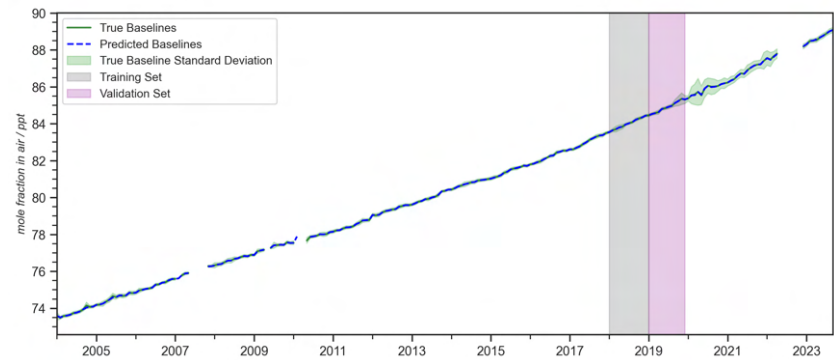


4.4.2 CF₄

30 Mole fraction time series

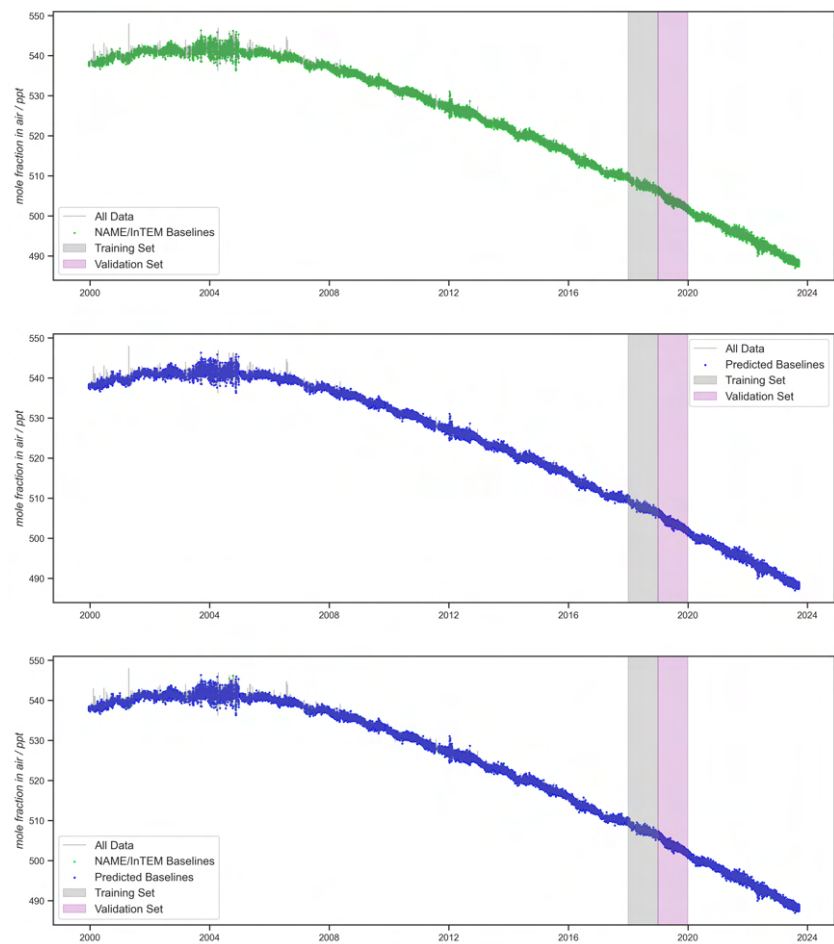


Monthly means

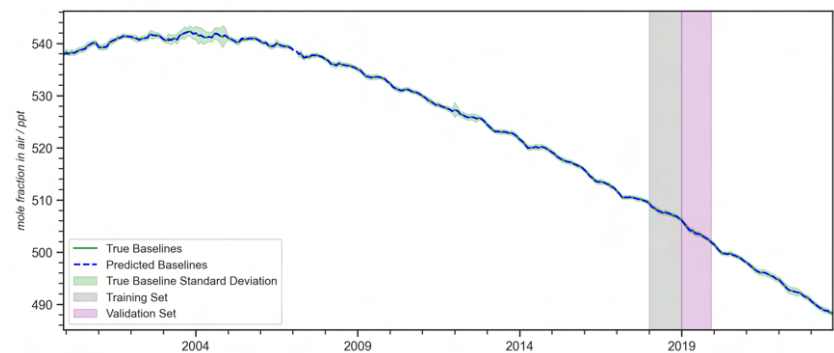


4.4.3 CFC-12

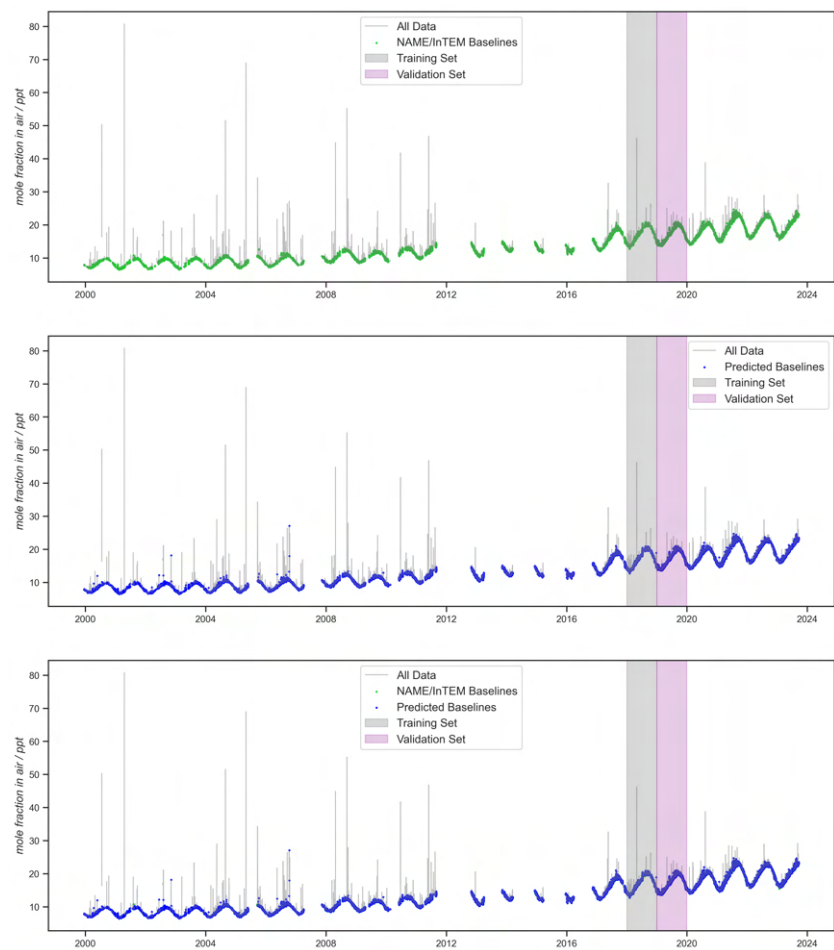
Mole fraction time series



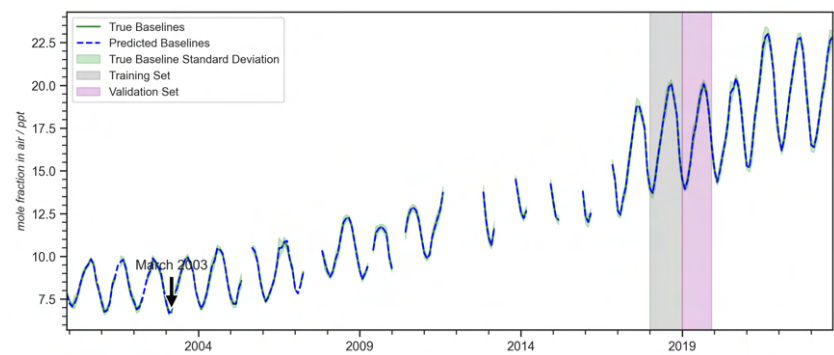
Monthly means



Mole fraction time series

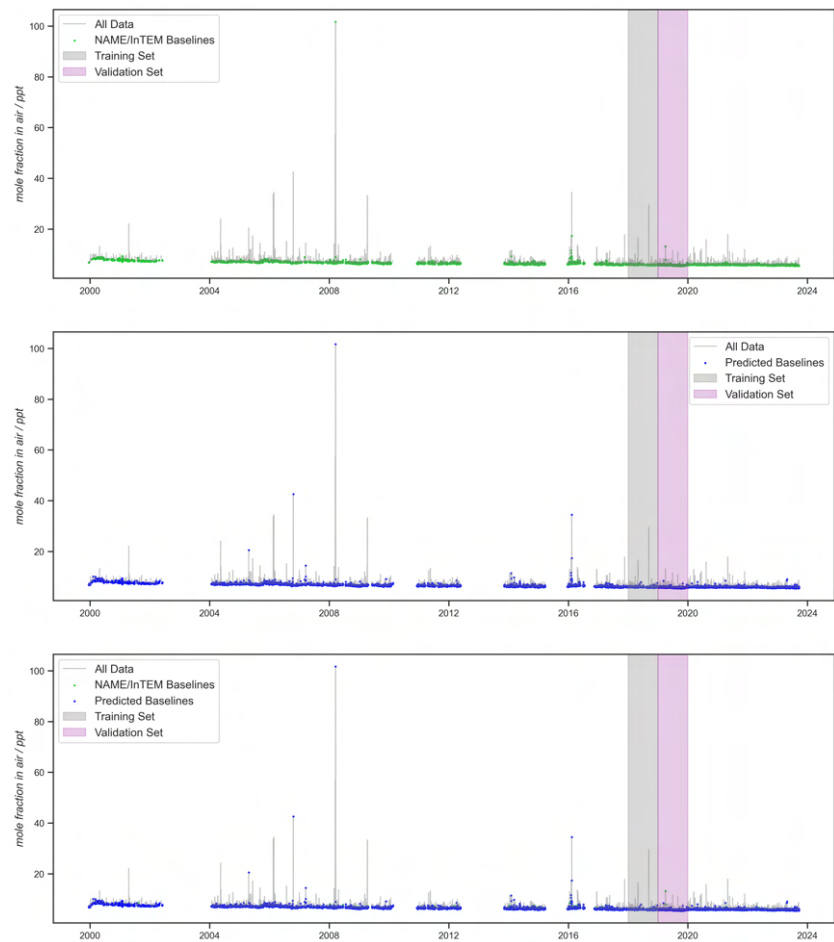


Monthly means

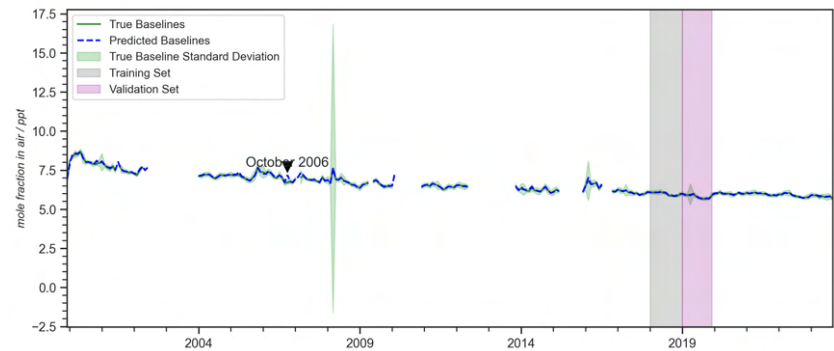


4.4.5 CH₃Br

Mole fraction time series

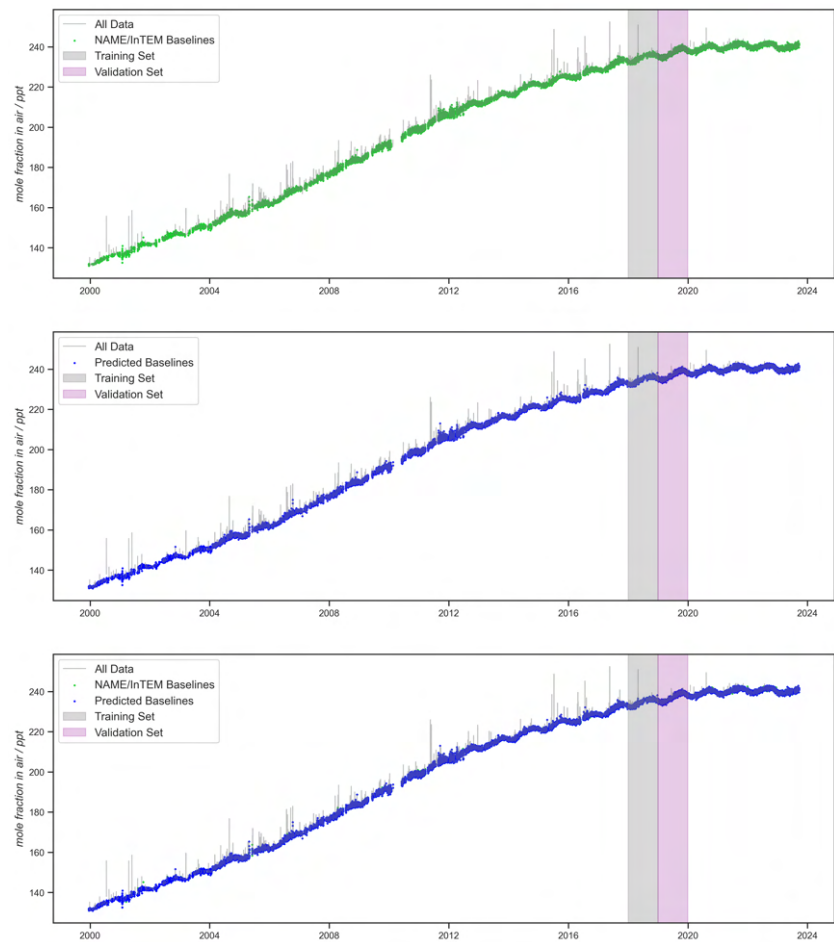


40 Monthly means

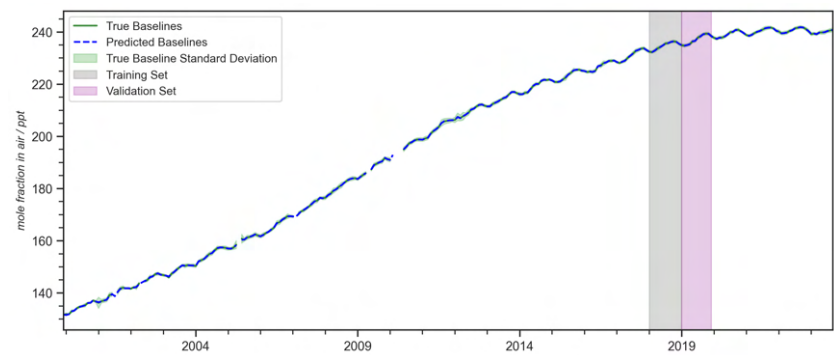


4.4.6 HCFC-22

Mole fraction time series

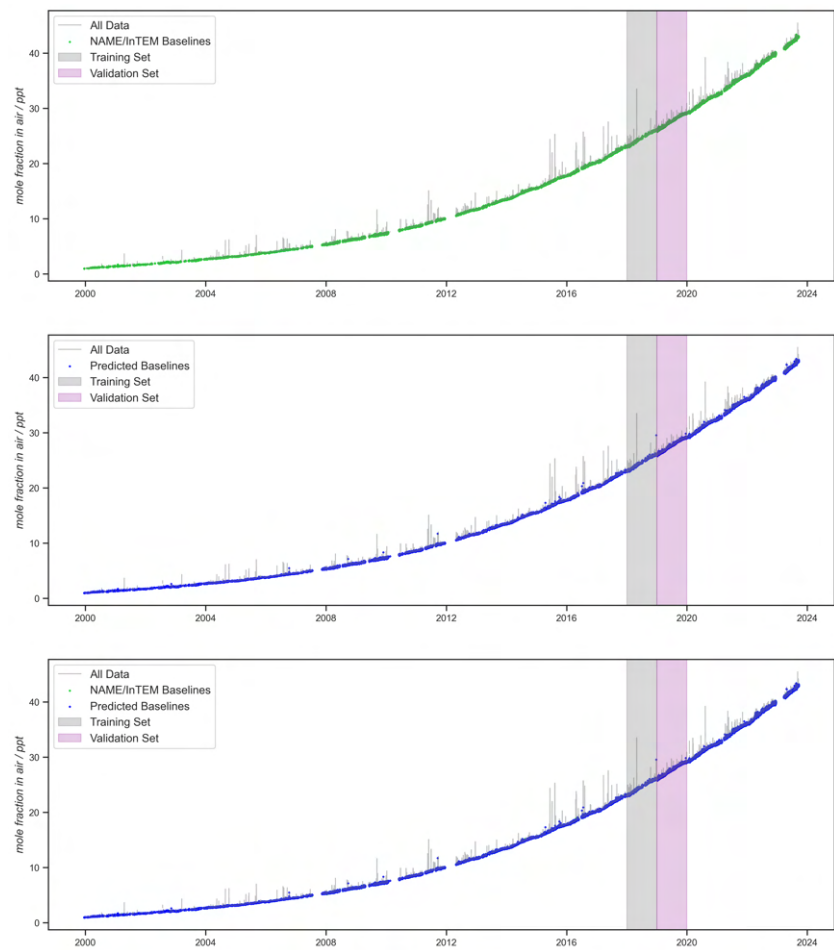


Monthly means

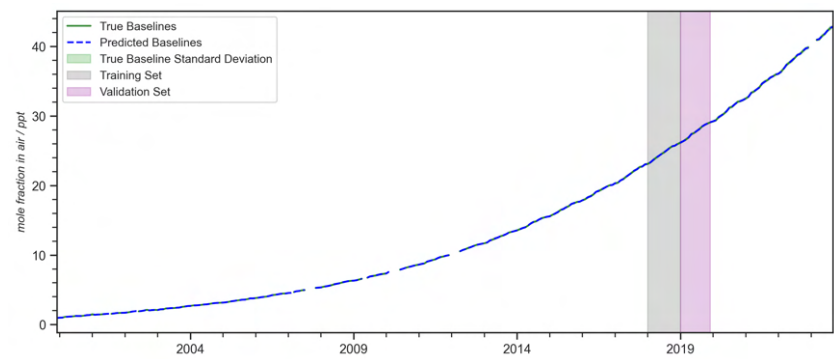


4.4.7 HFC-125

45 Mole fraction time series

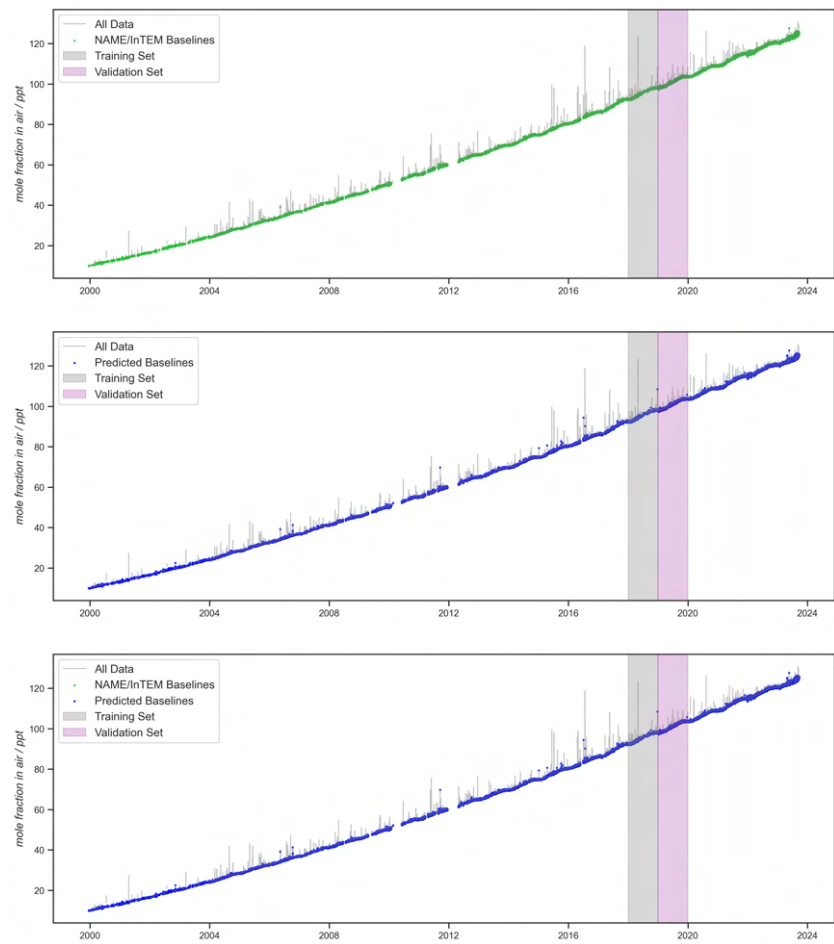


Monthly means

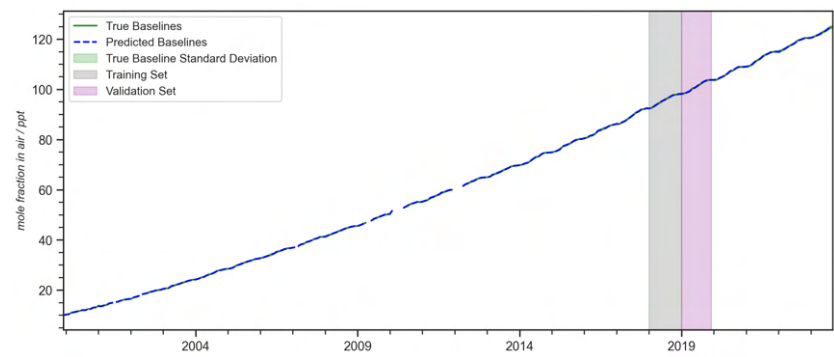


4.4.8 HFC-134a

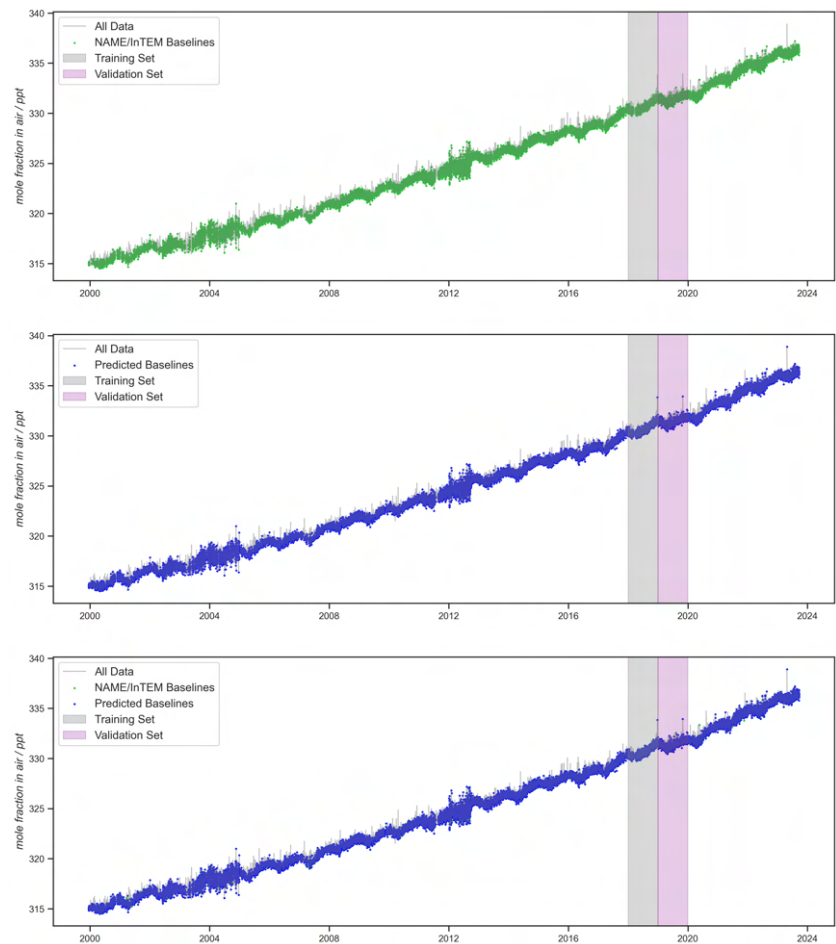
Mole fraction time series



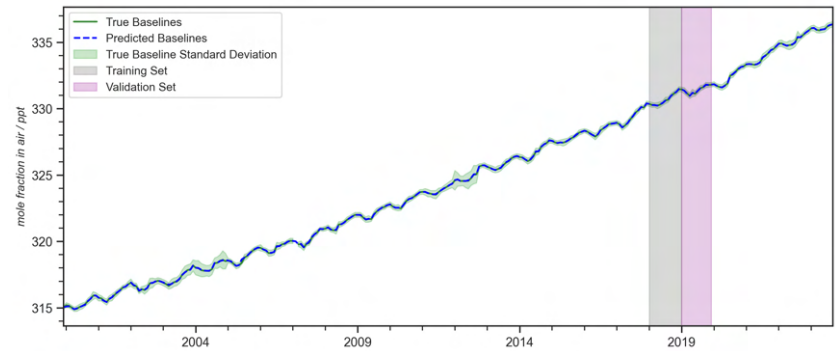
Monthly means



Mole fraction time series

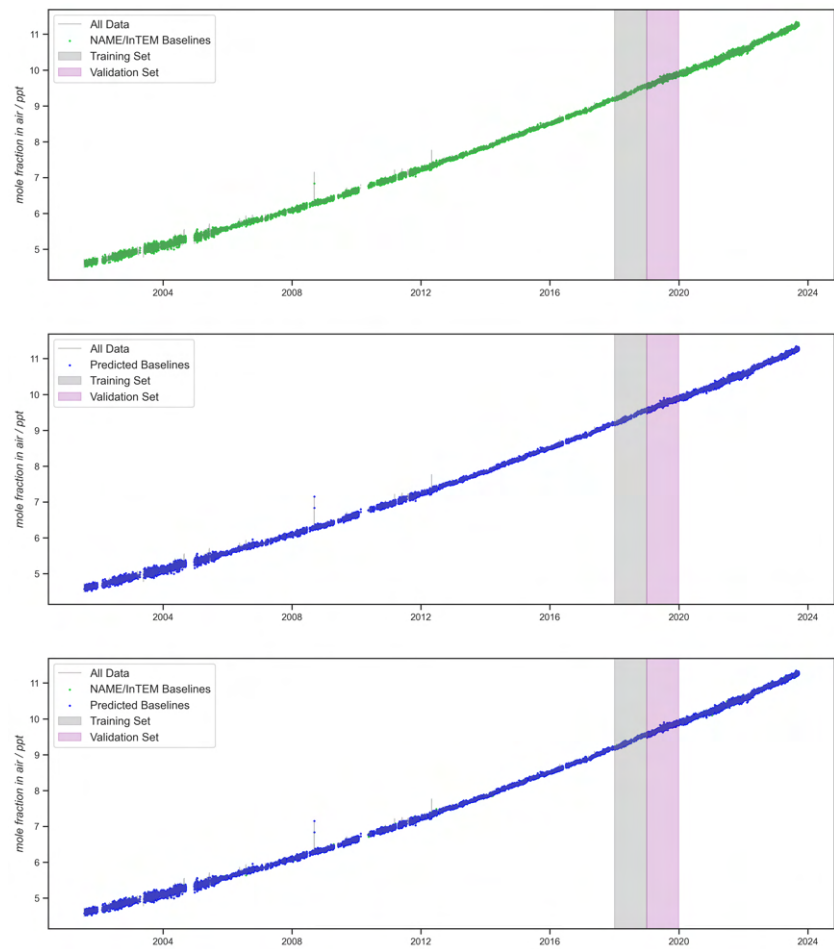


Monthly means

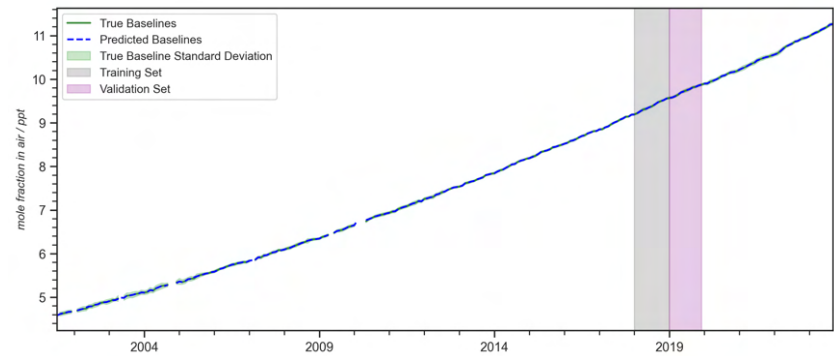


4.4.10 SF₆

Mole fraction time series



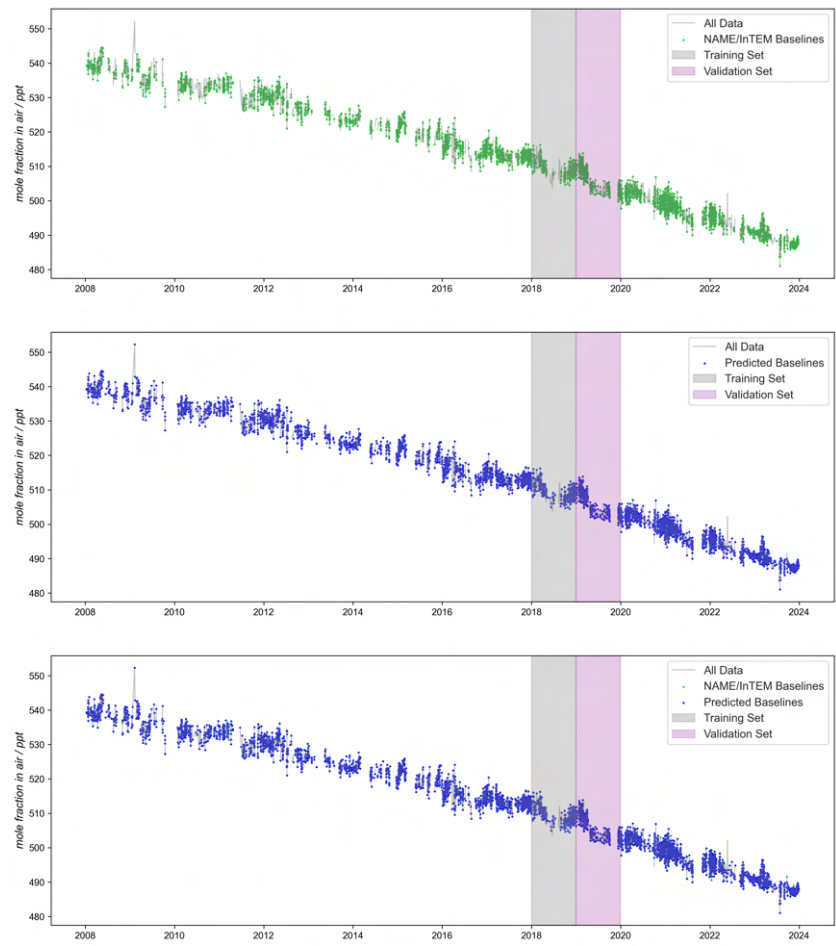
55 Monthly means



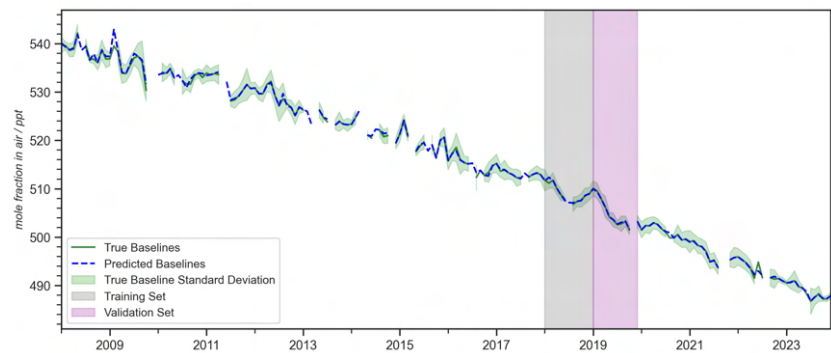
4.5 Monte Cimone, Italy

4.5.1 CFC-12

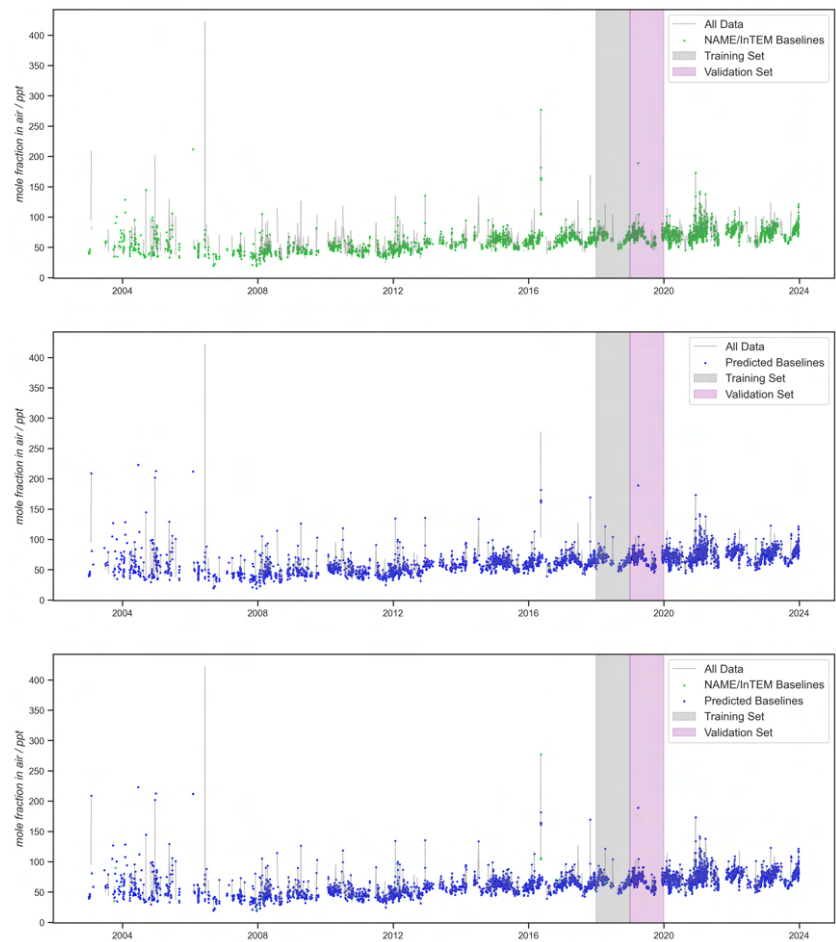
Mole fraction time series



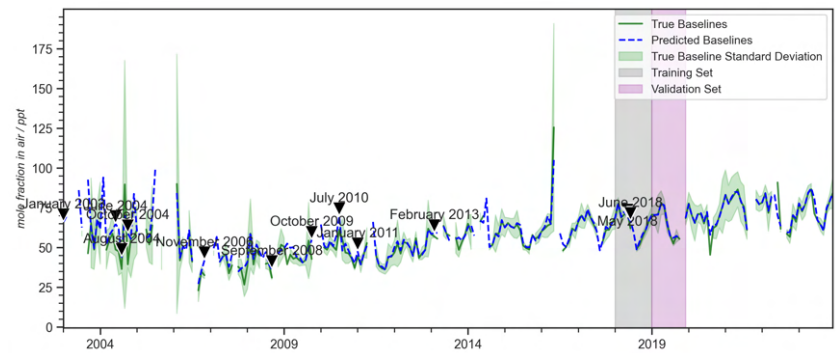
Monthly means



Mole fraction time series

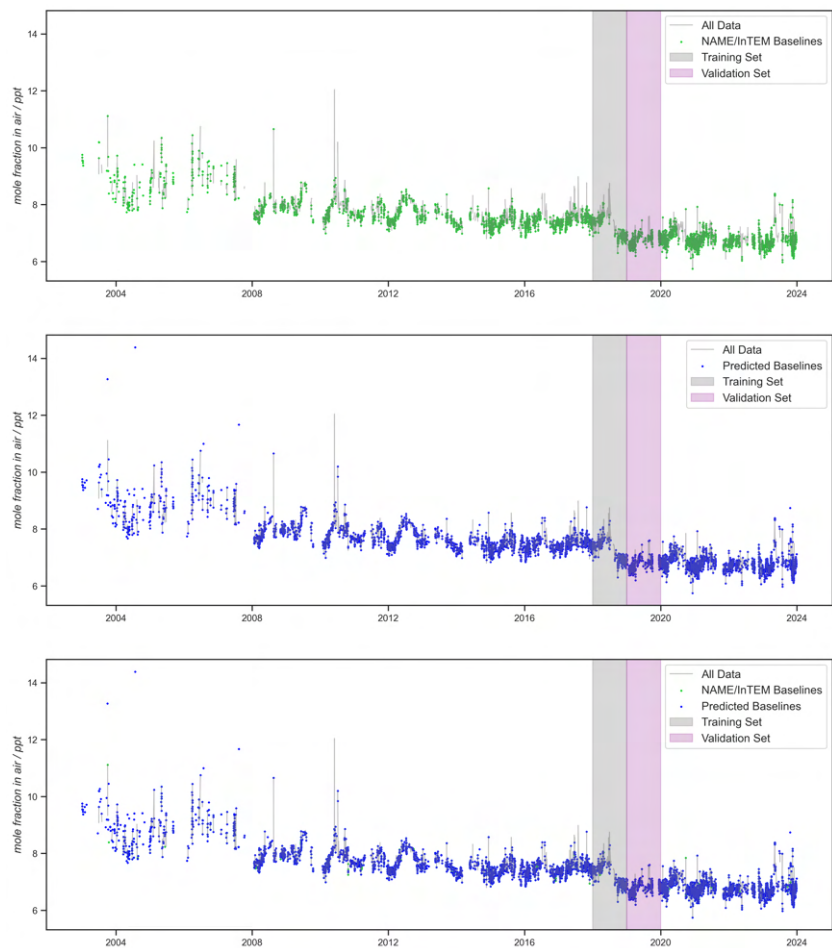


Monthly means

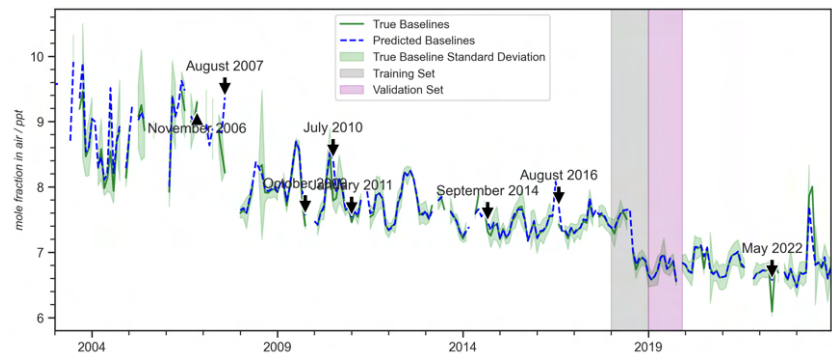


4.5.3 CH₃Br

Mole fraction time series

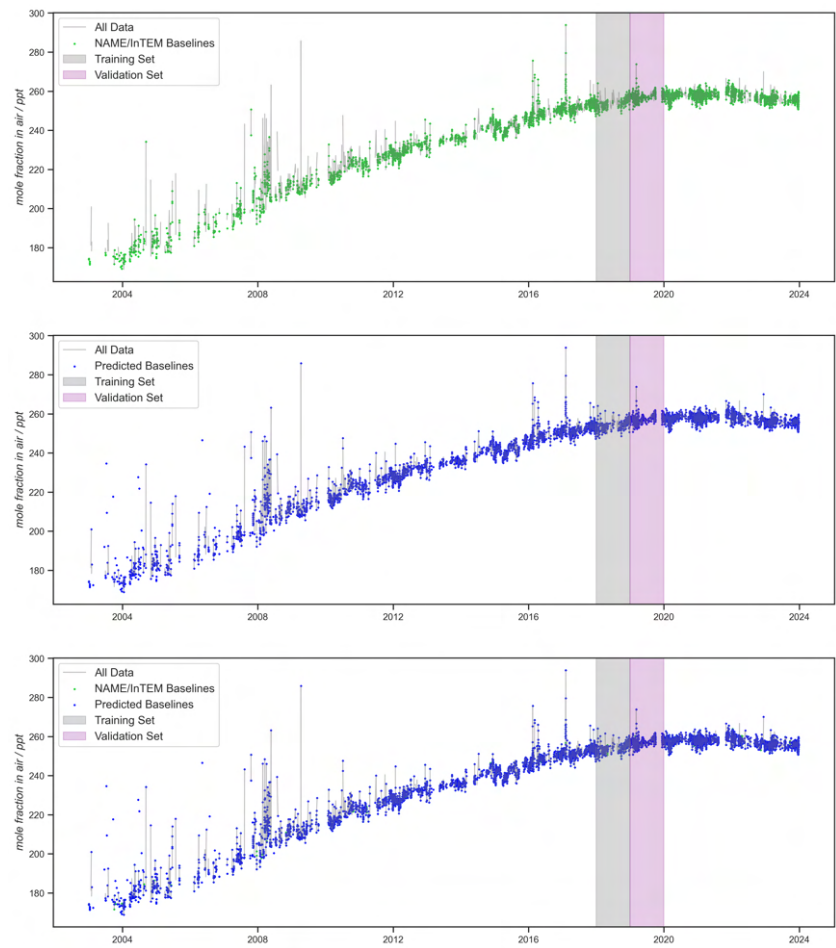


65 Monthly means

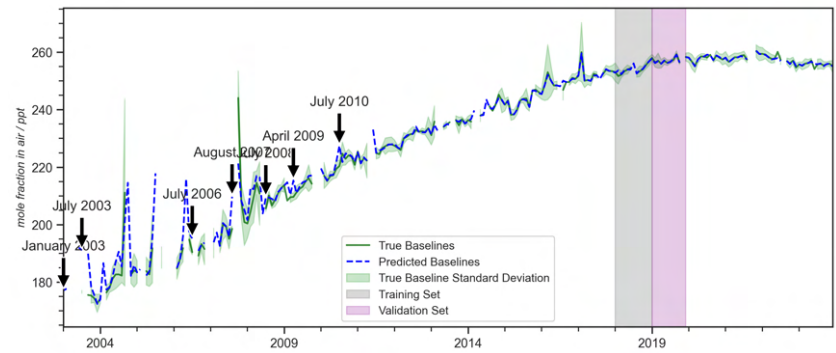


4.5.4 HCFC-22

Mole fraction time series

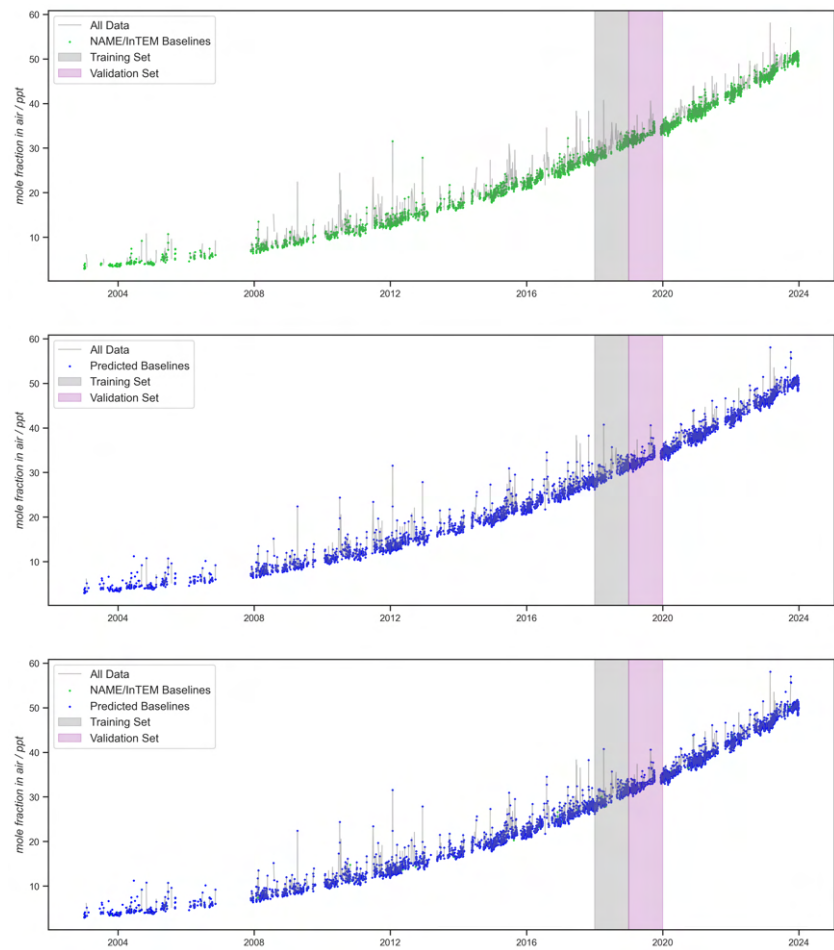


Monthly means

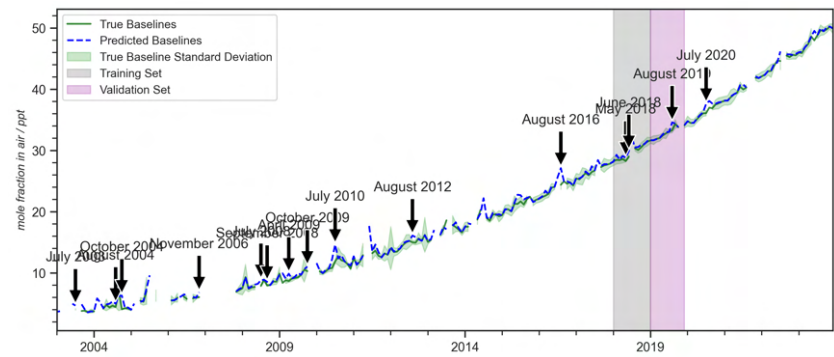


4.5.5 HFC-125

70 Mole fraction time series

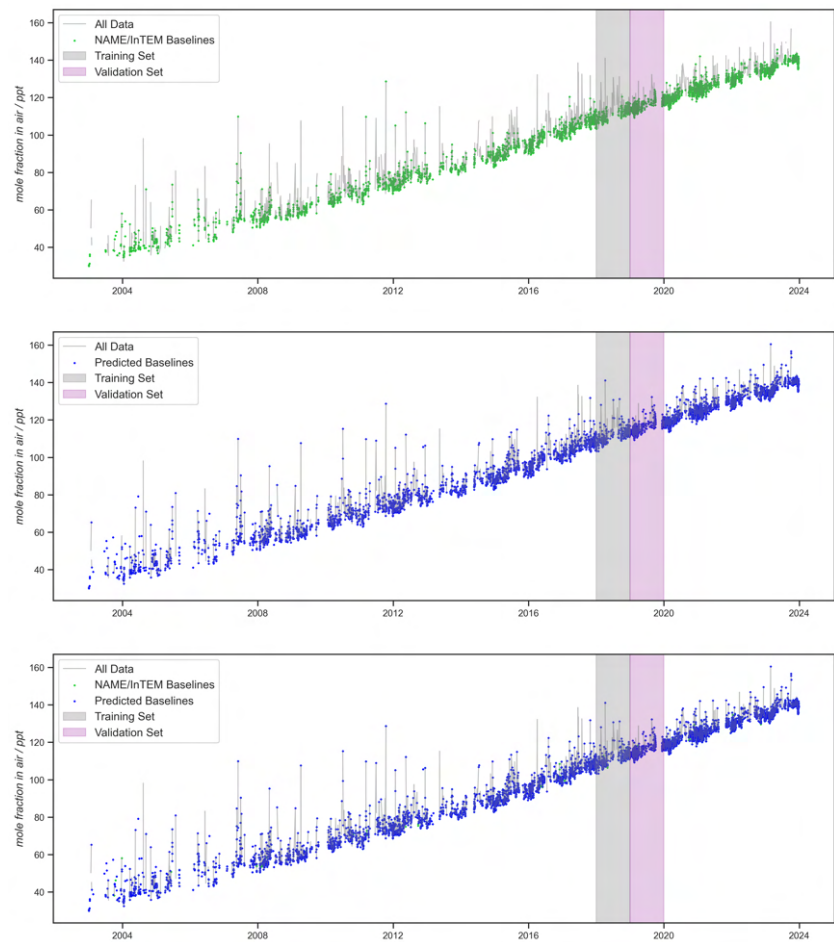


Monthly means

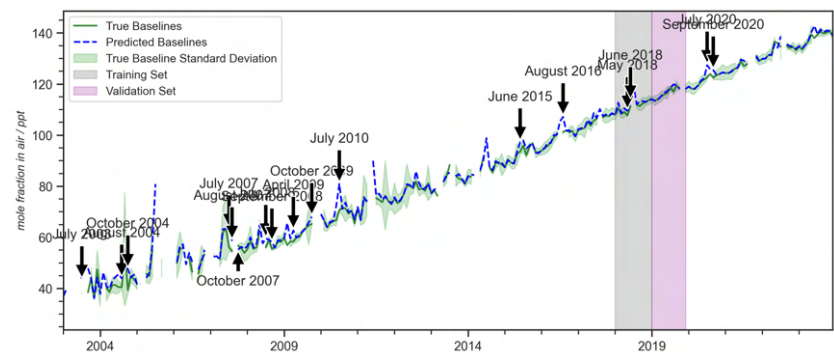


4.5.6 HFC-134a

Mole fraction time series

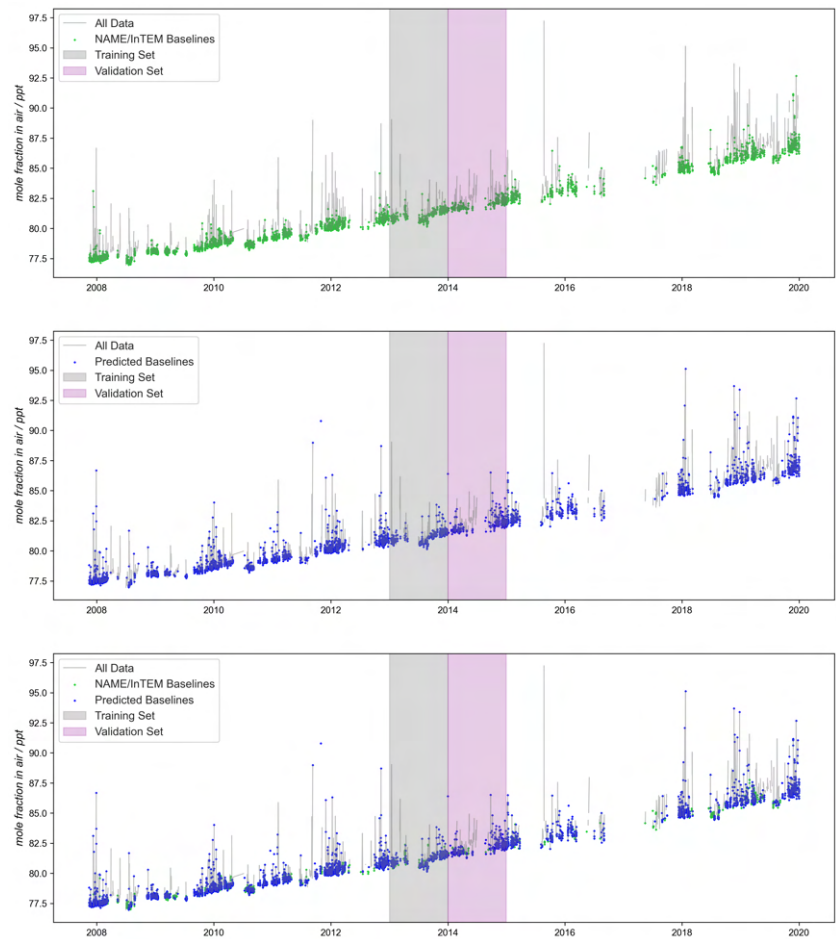


Monthly means

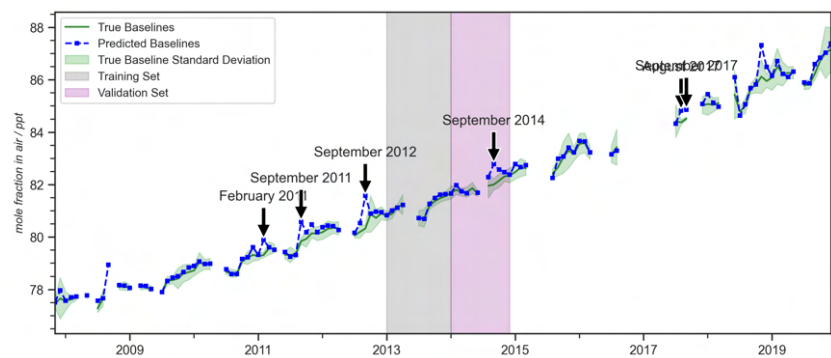


4.6.1 CF₄

Mole fraction time series

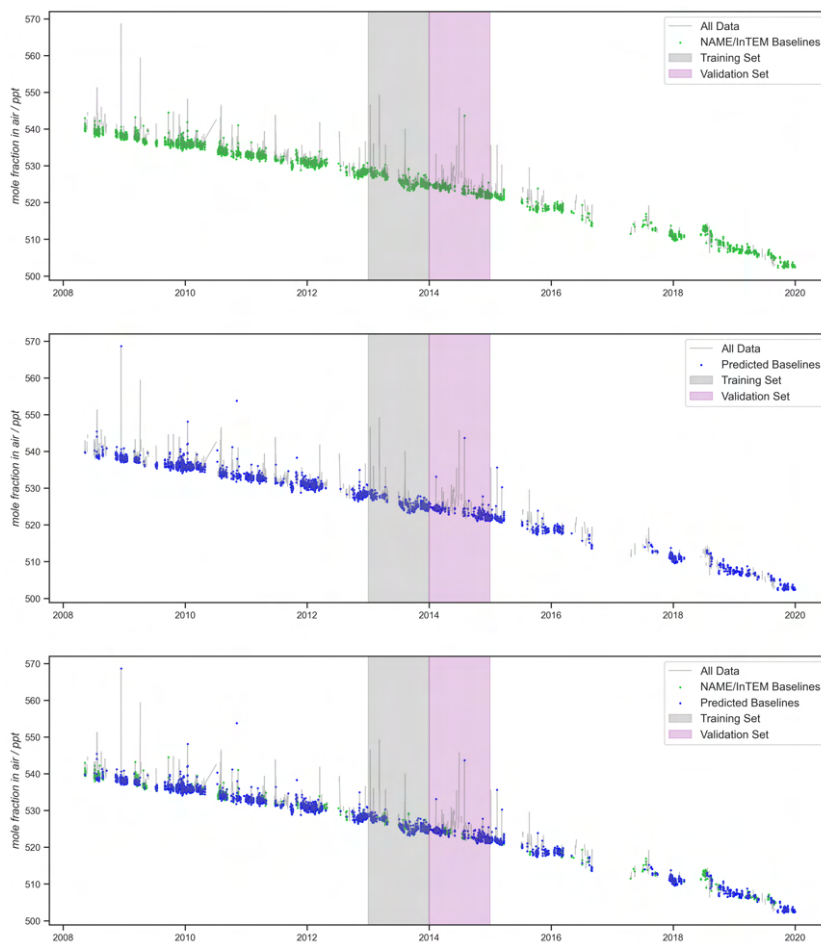


Monthly means

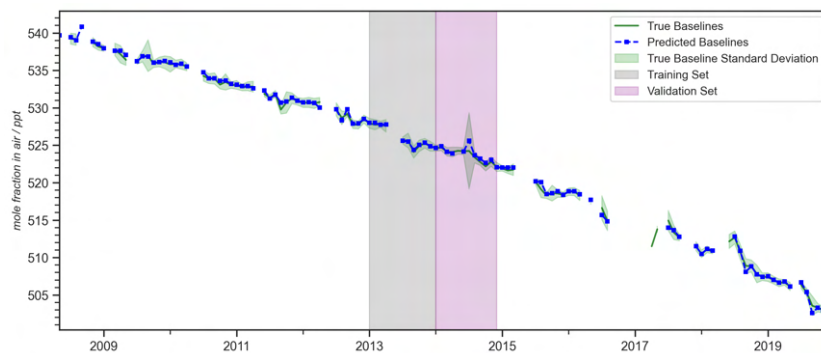


4.6.2 CFC-12

80 Mole fraction time series

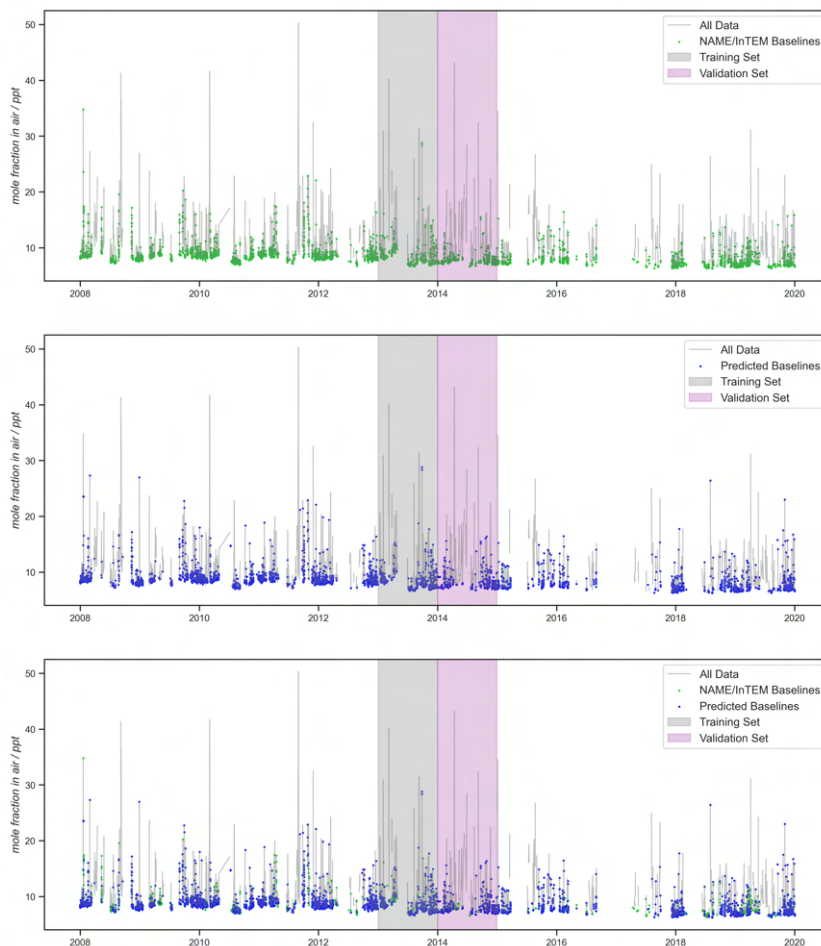


Monthly means

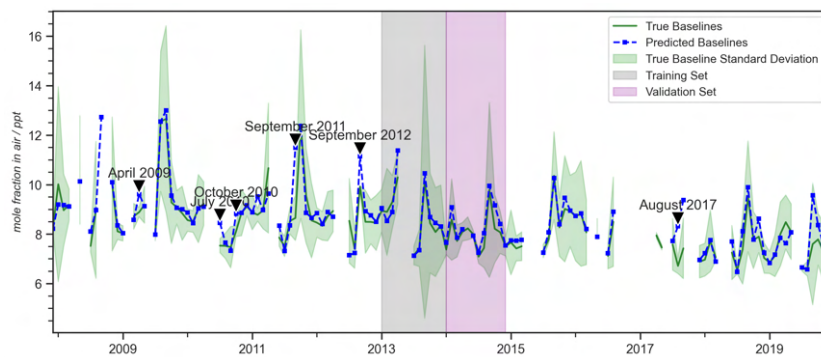


4.6.3 CH₃Br

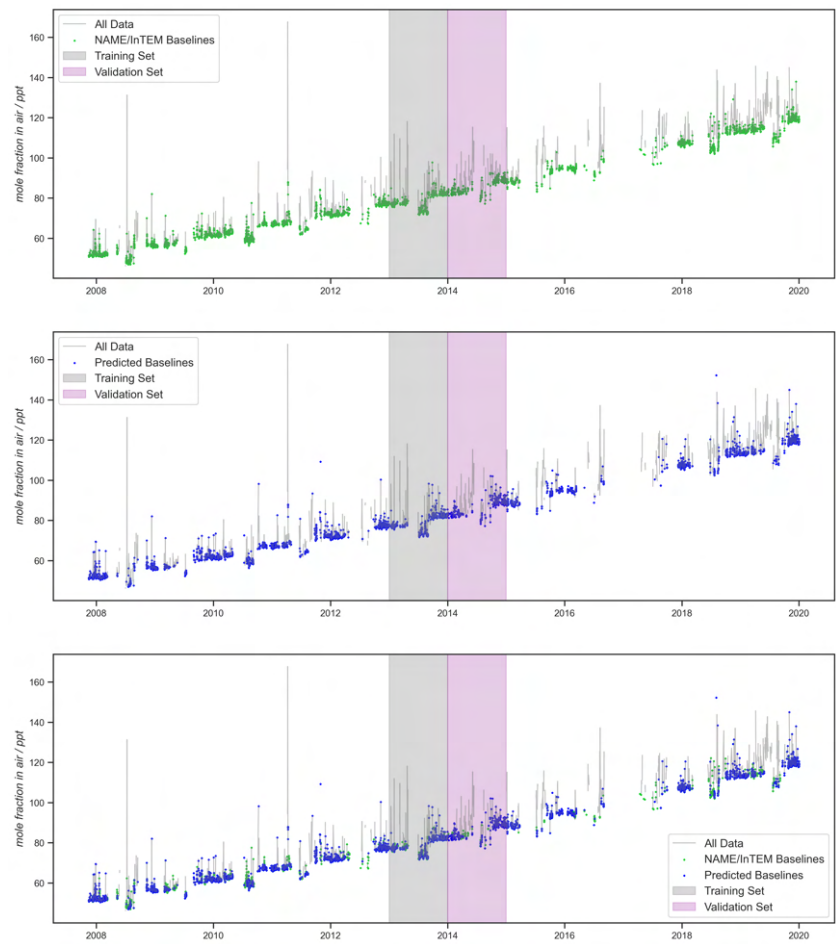
Mole fraction time series



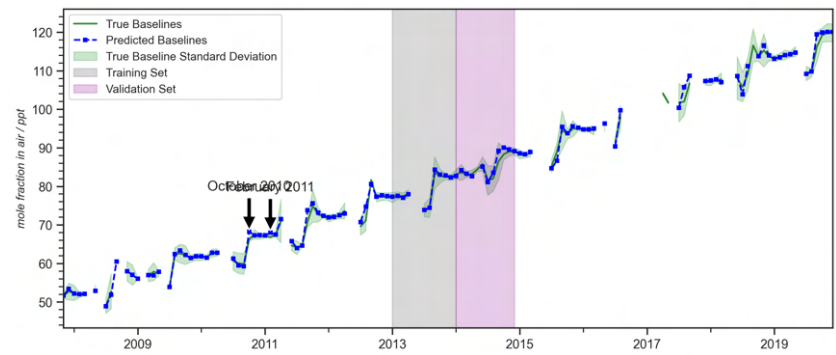
Monthly means



Mole fraction time series

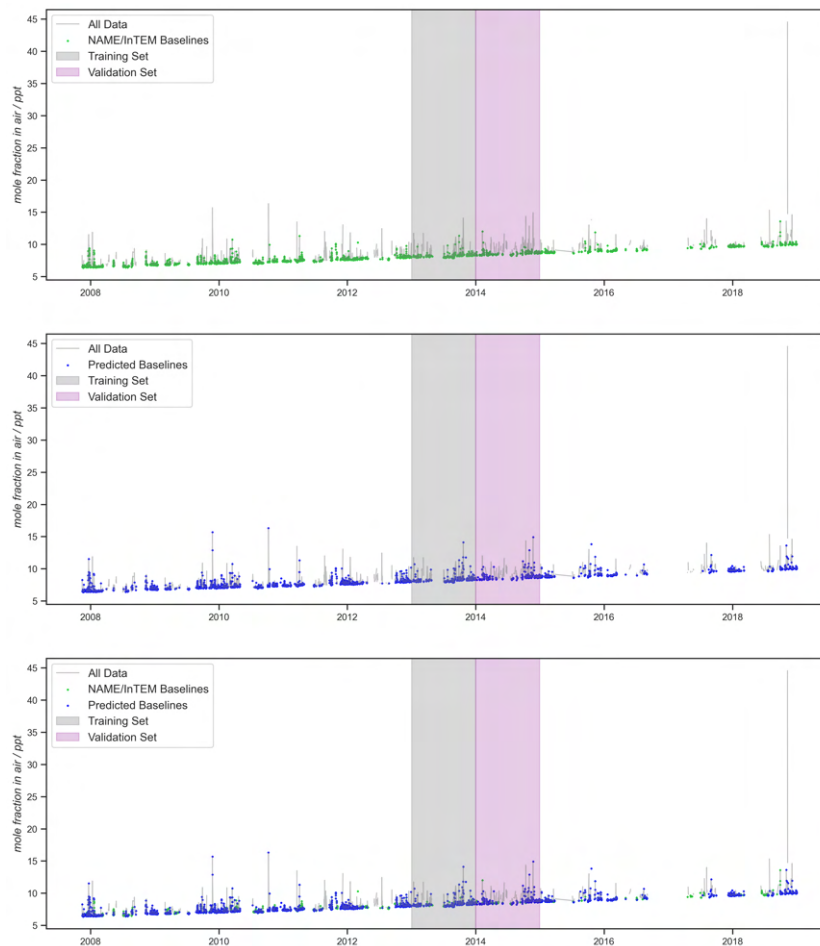


Monthly means

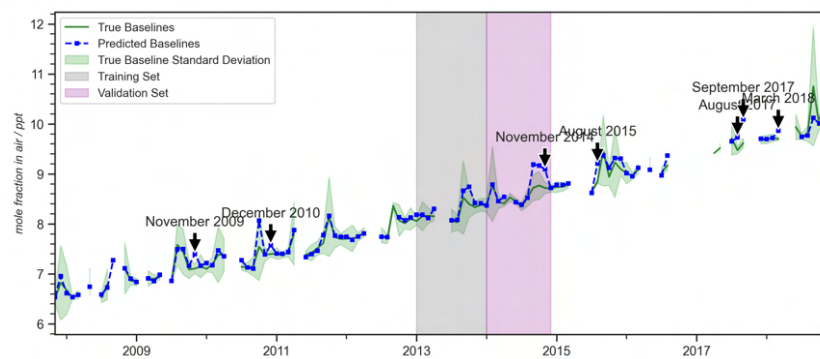


4.6.5 SF₆

Mole fraction time series



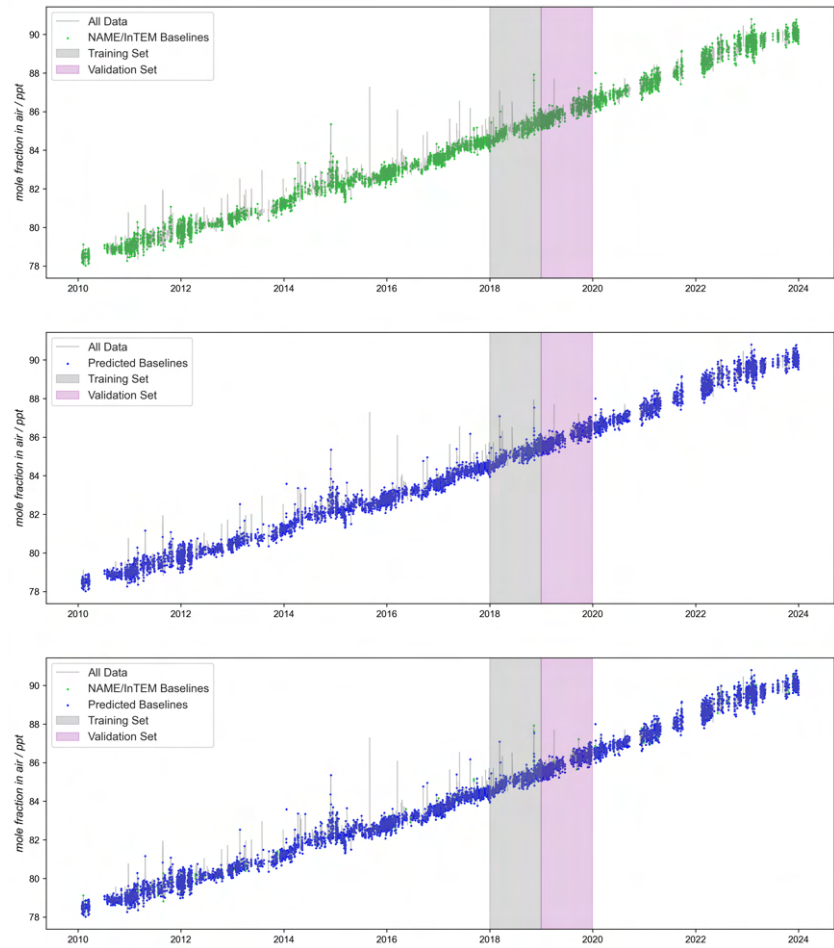
90 Monthly means



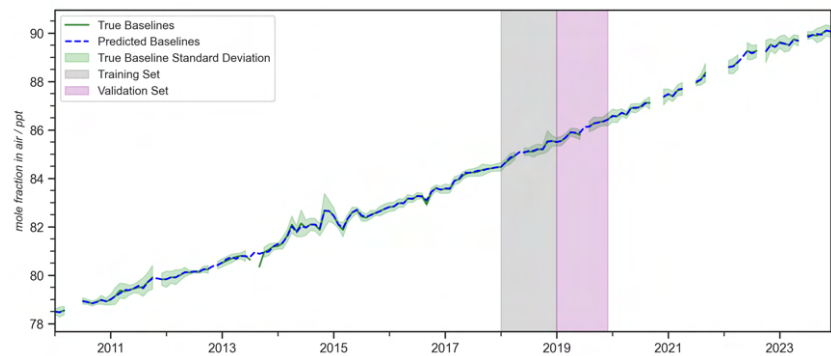
4.7 Jungfraujoch, Switzerland

4.7.1 CF₄

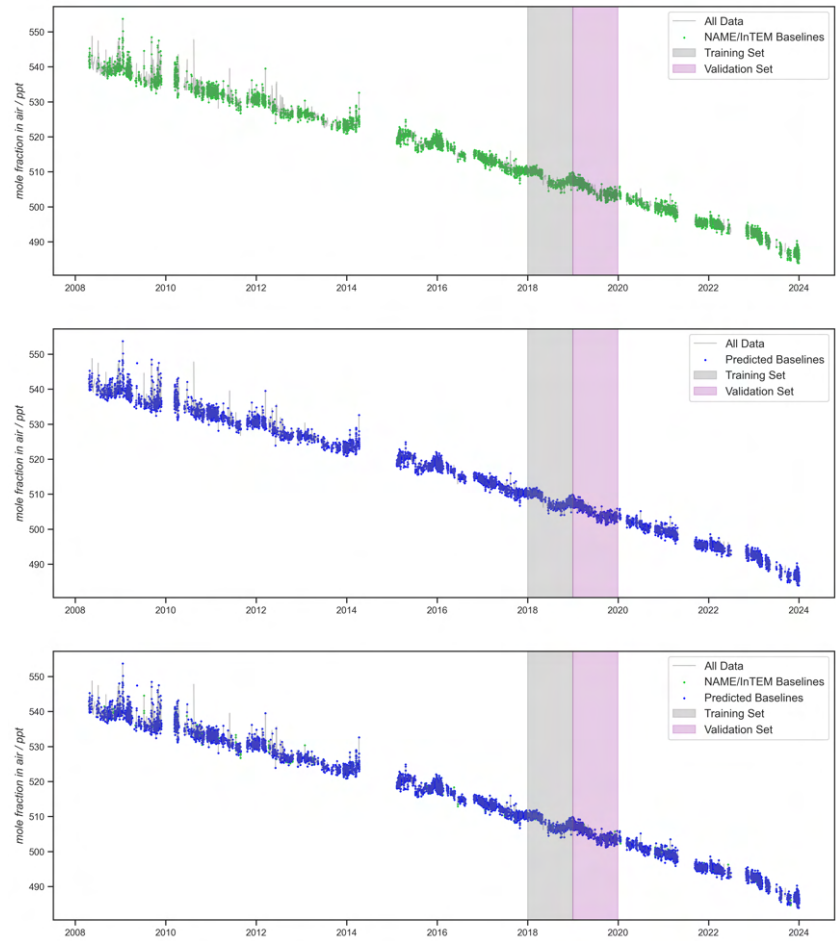
Mole fraction time series



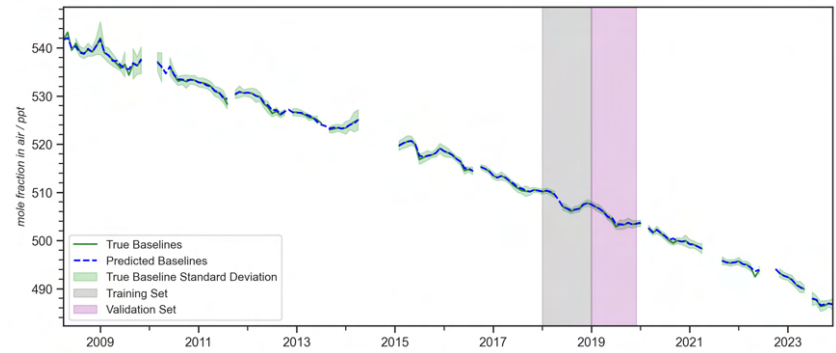
Monthly means



Mole fraction time series

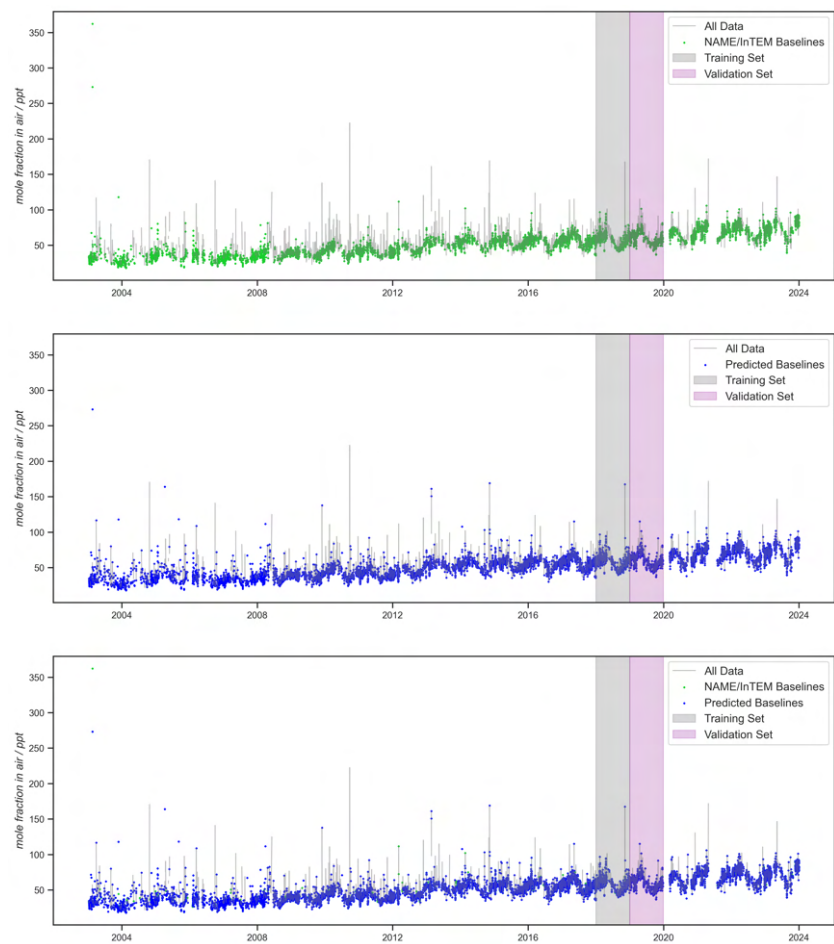


Monthly means

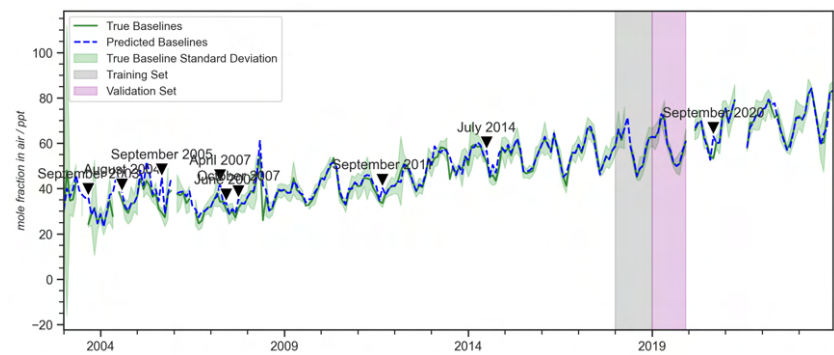


4.7.3 CH₂Cl₂

Mole fraction time series

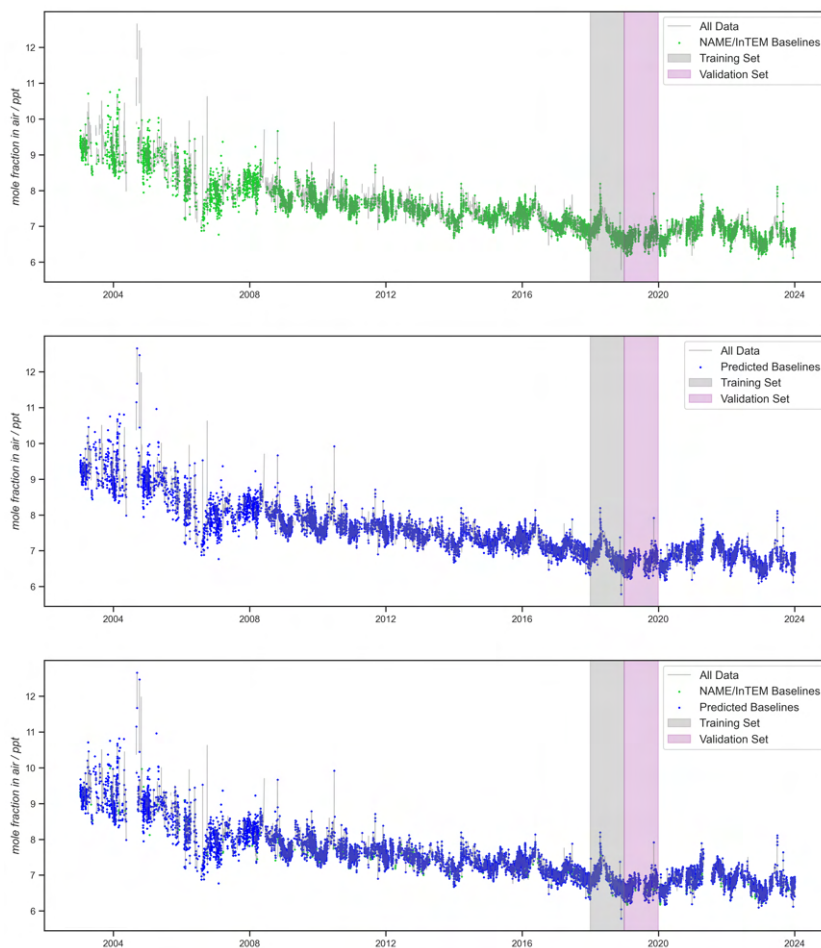


100 Monthly means

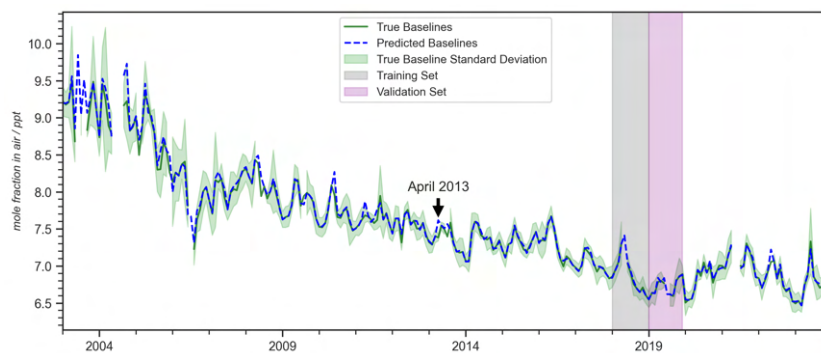


4.7.4 CH₃Br

Mole fraction time series

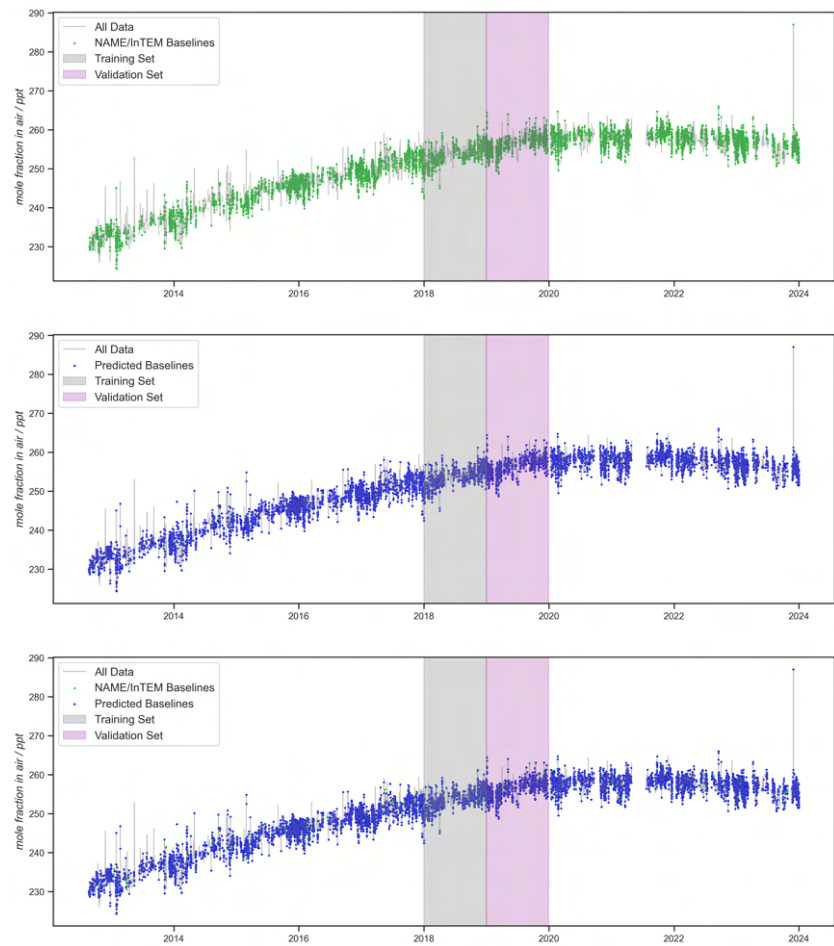


Monthly means

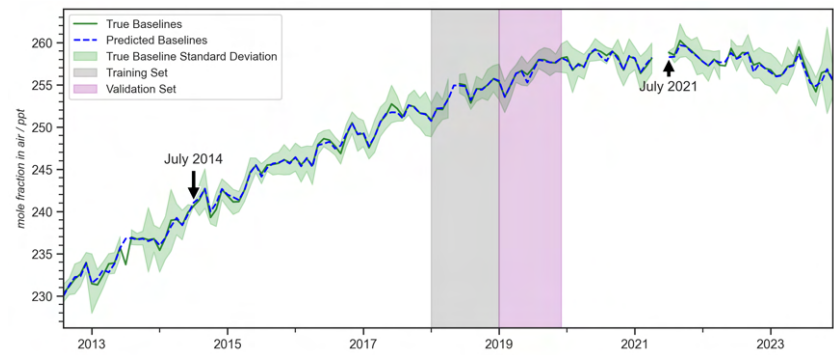


4.7.5 HCFC-22

105 Mole fraction time series

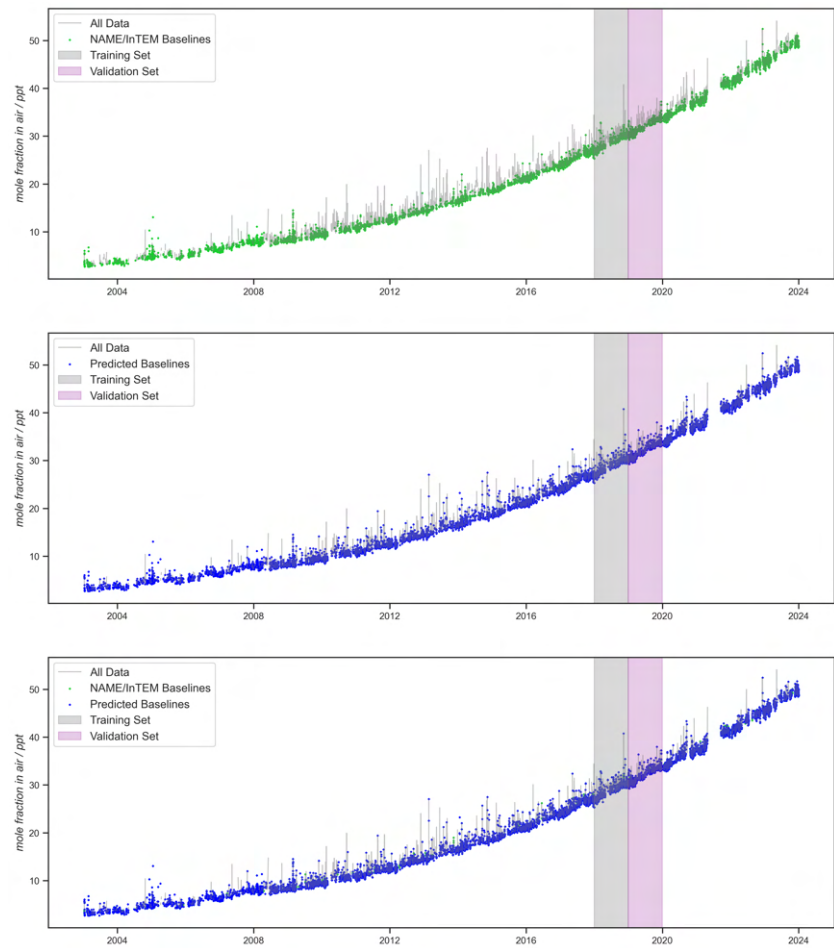


Monthly means

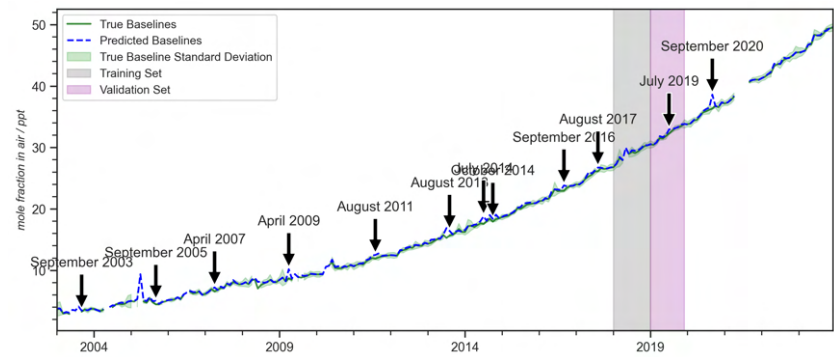


4.7.6 HFC-125

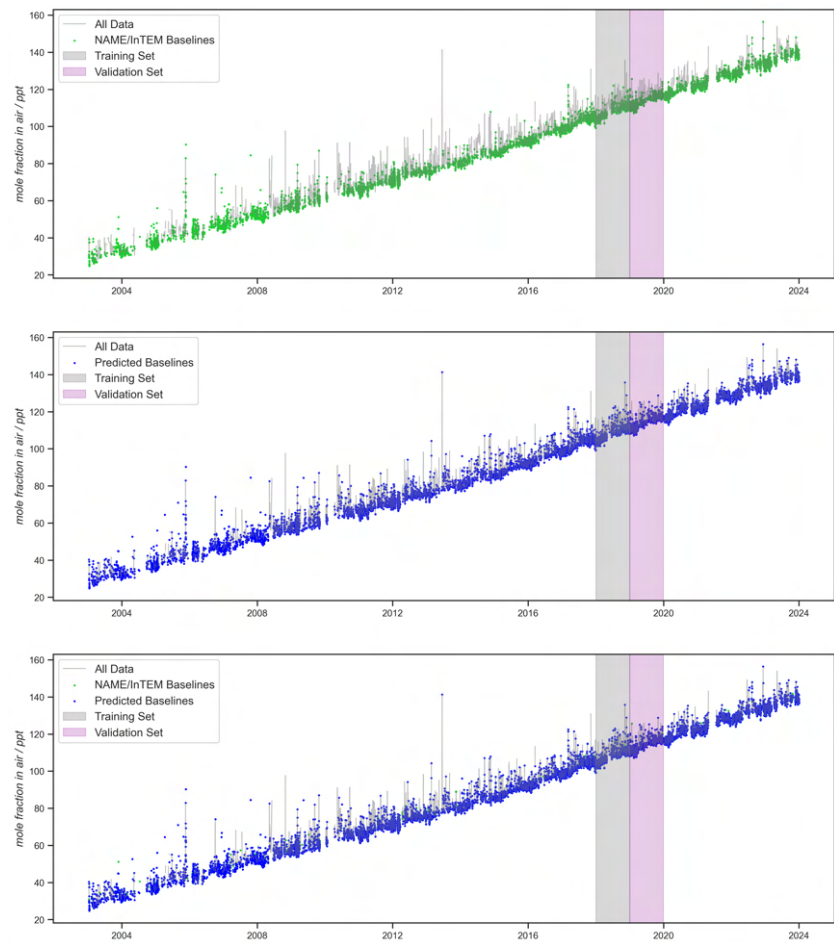
Mole fraction time series



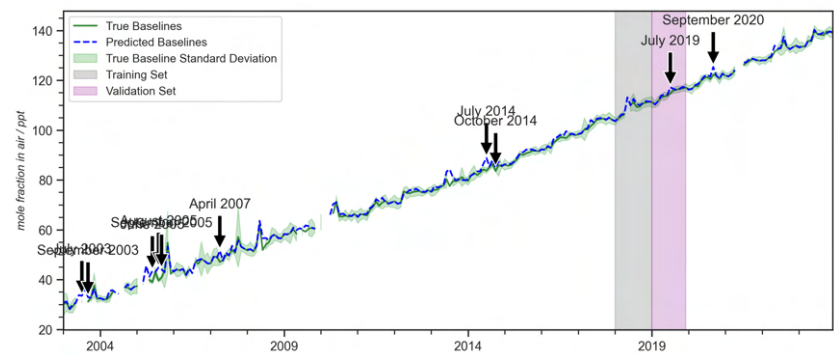
Monthly means



Mole fraction time series

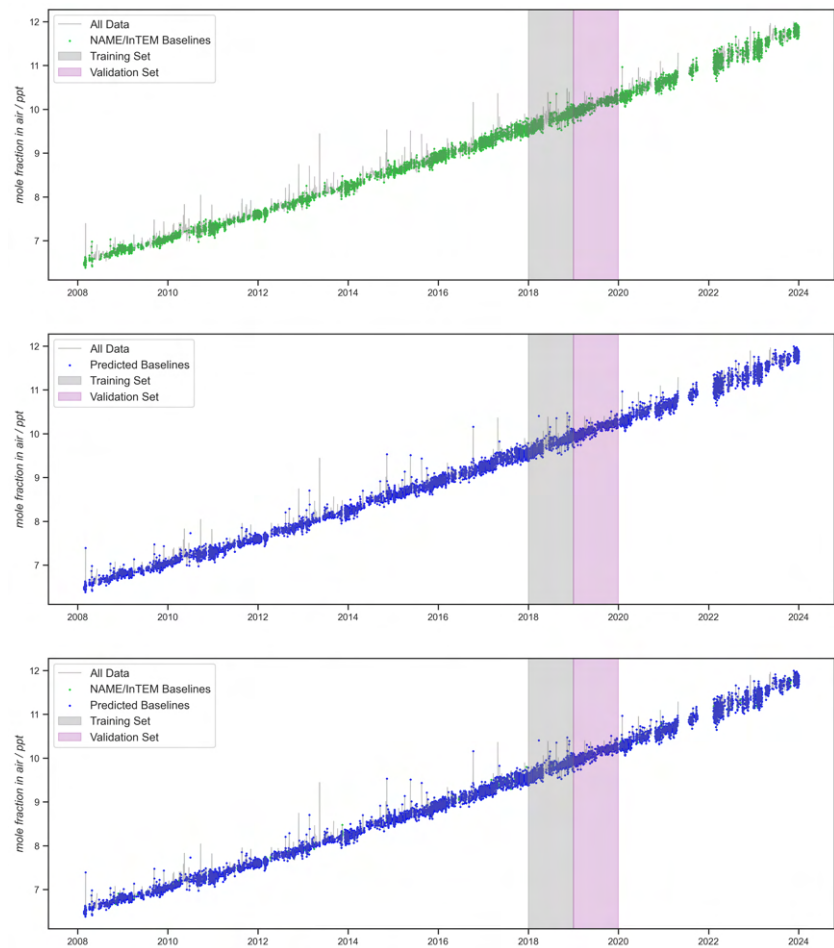


Monthly means

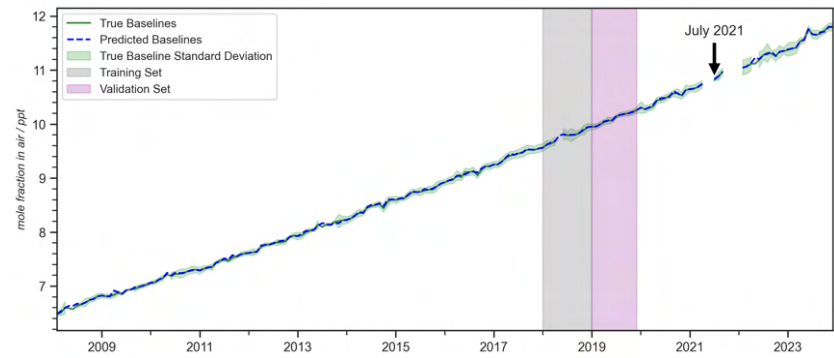


4.7.8 SF₆

Mole fraction time series



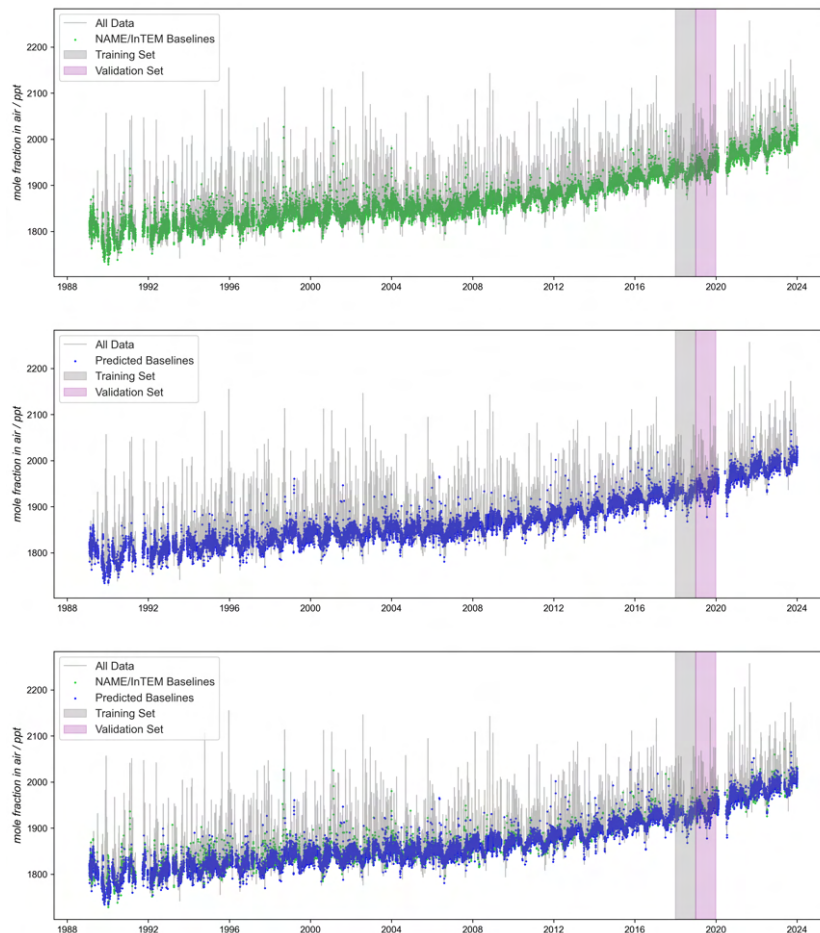
115 Monthly means



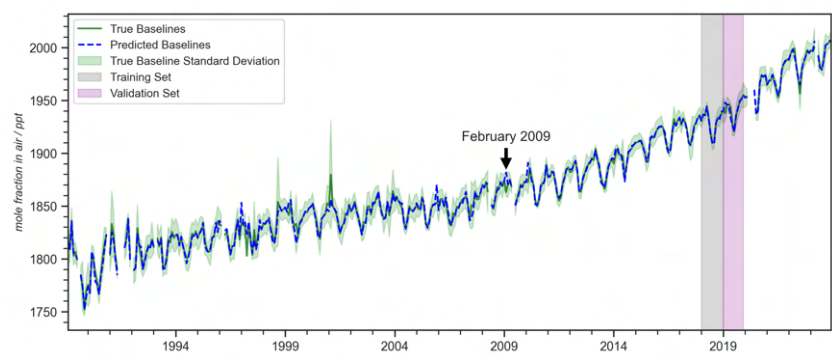
4.8 Mace Head, Ireland

4.8.1 CH₄

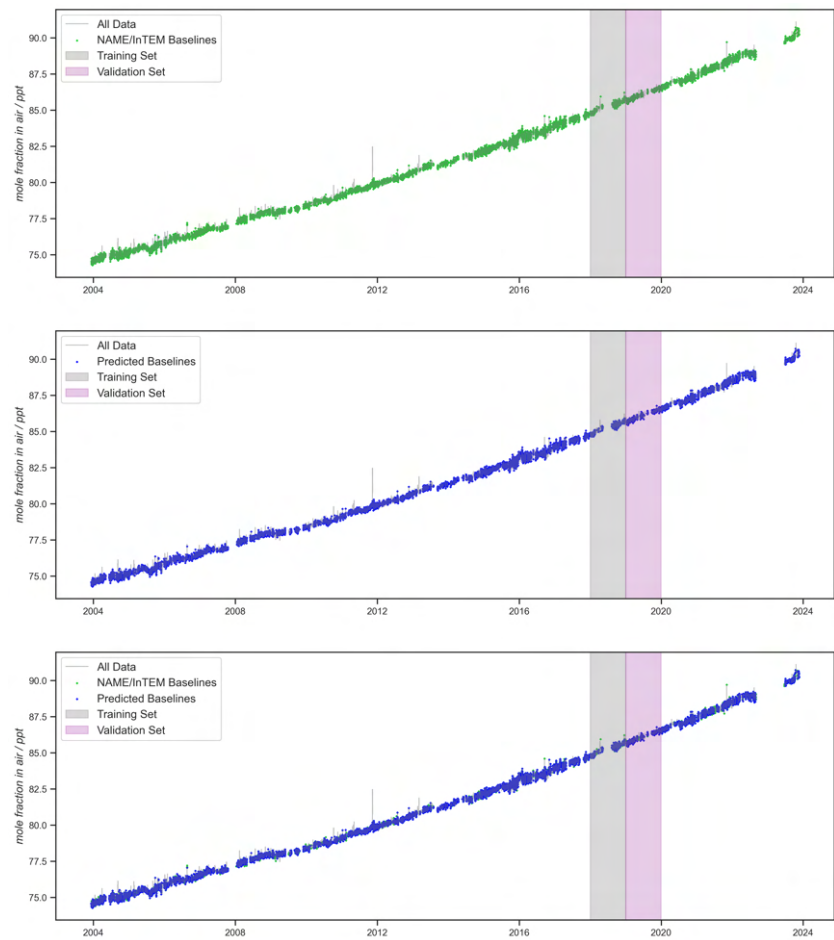
Mole fraction time series



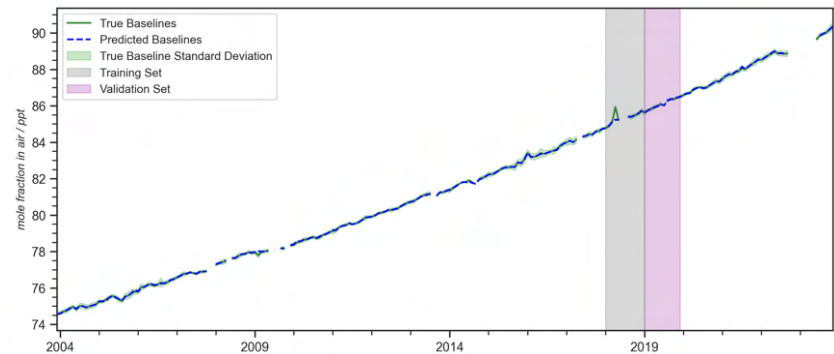
Monthly means



Mole fraction time series

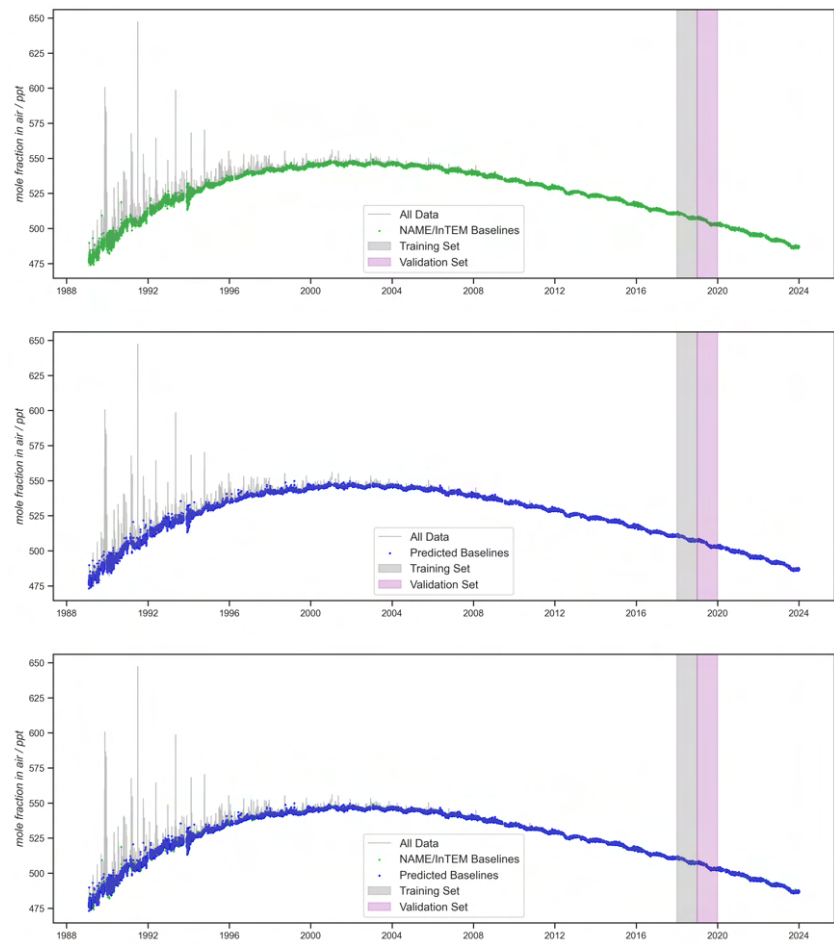


Monthly means

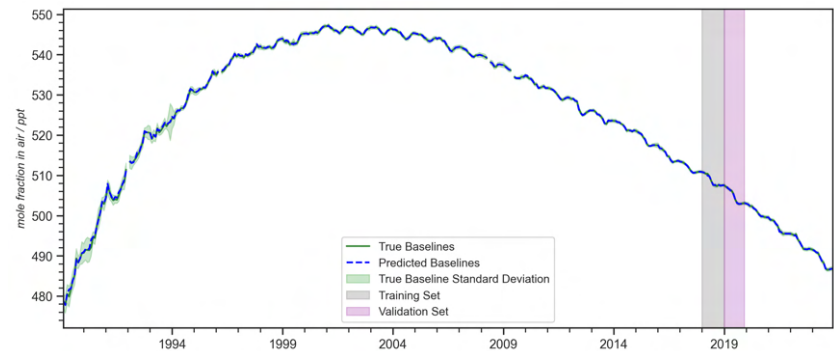


4.8.3 CFC-12

Mole fraction time series

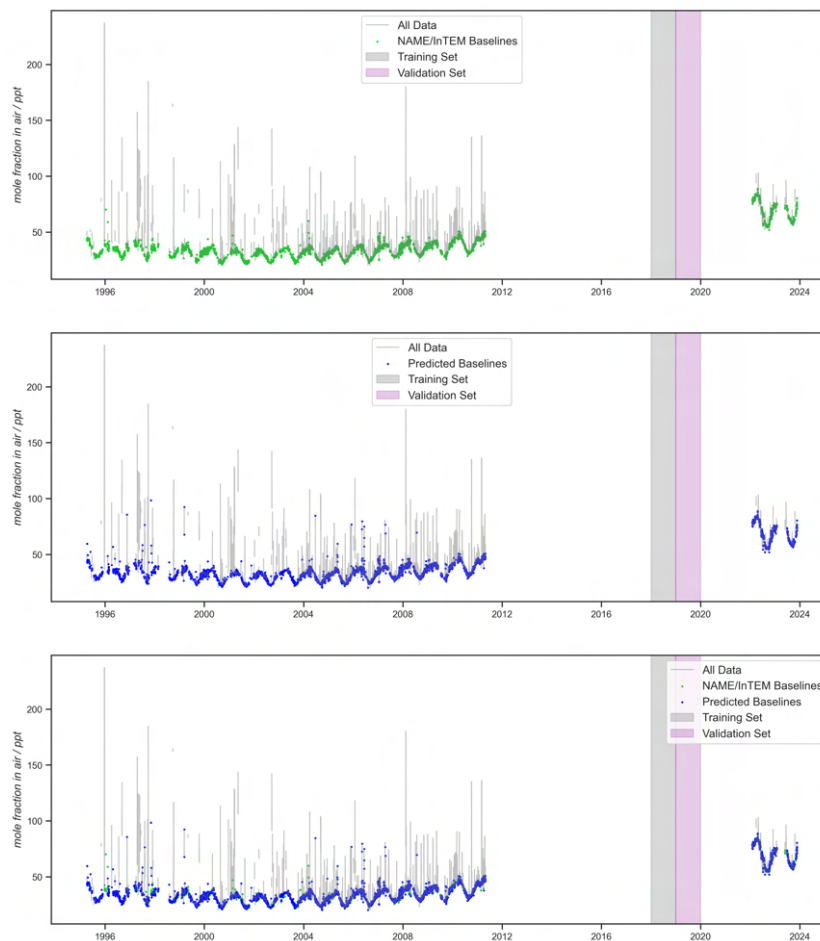


125 Monthly means

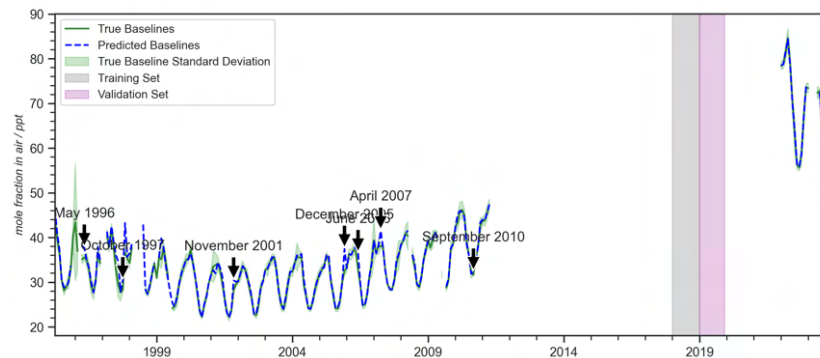


4.8.4 CH₂Cl₂

Mole fraction time series

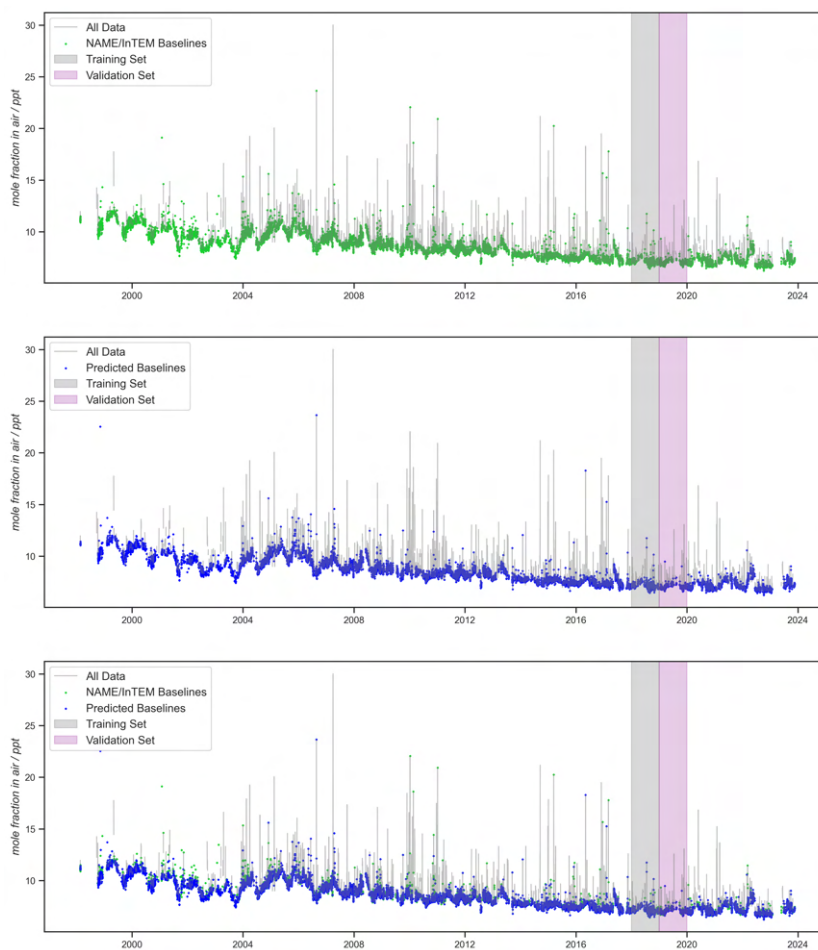


Monthly means

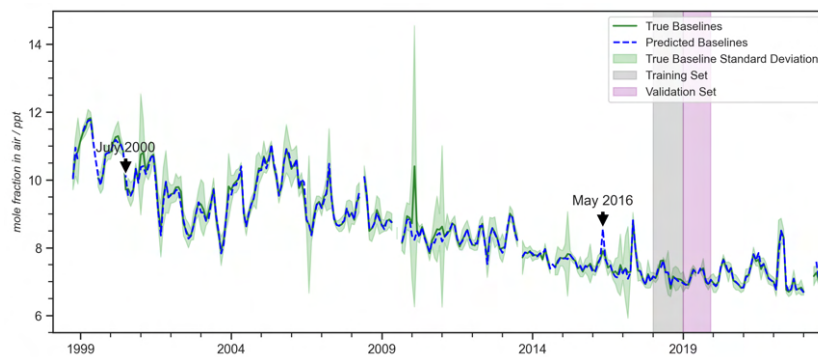


4.8.5 CH₃Br

130 Mole fraction time series

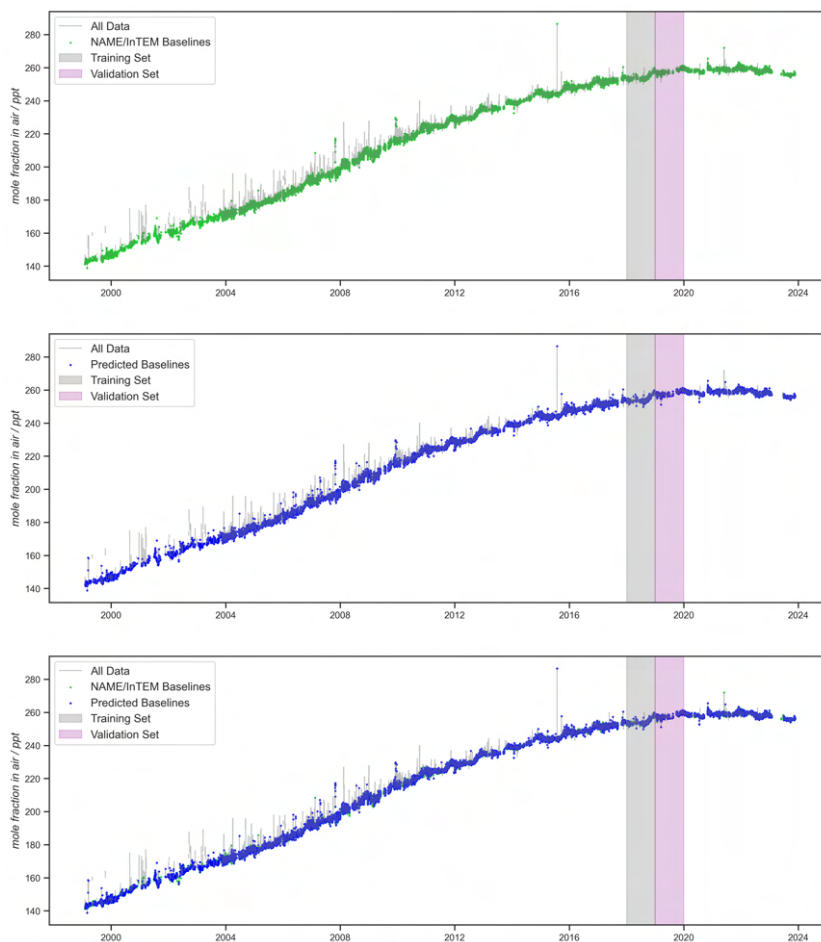


Monthly means

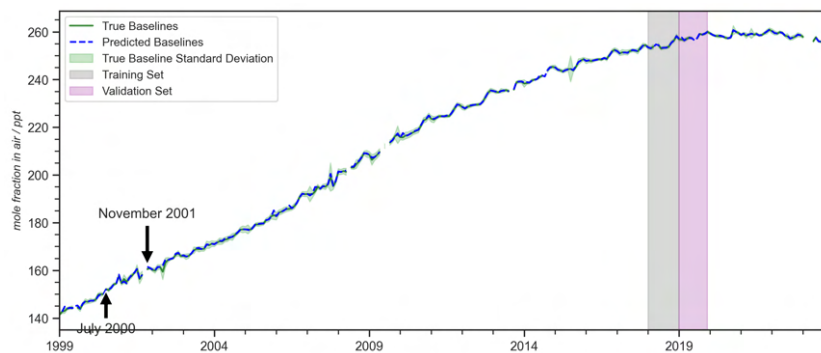


4.8.6 HCFC-22

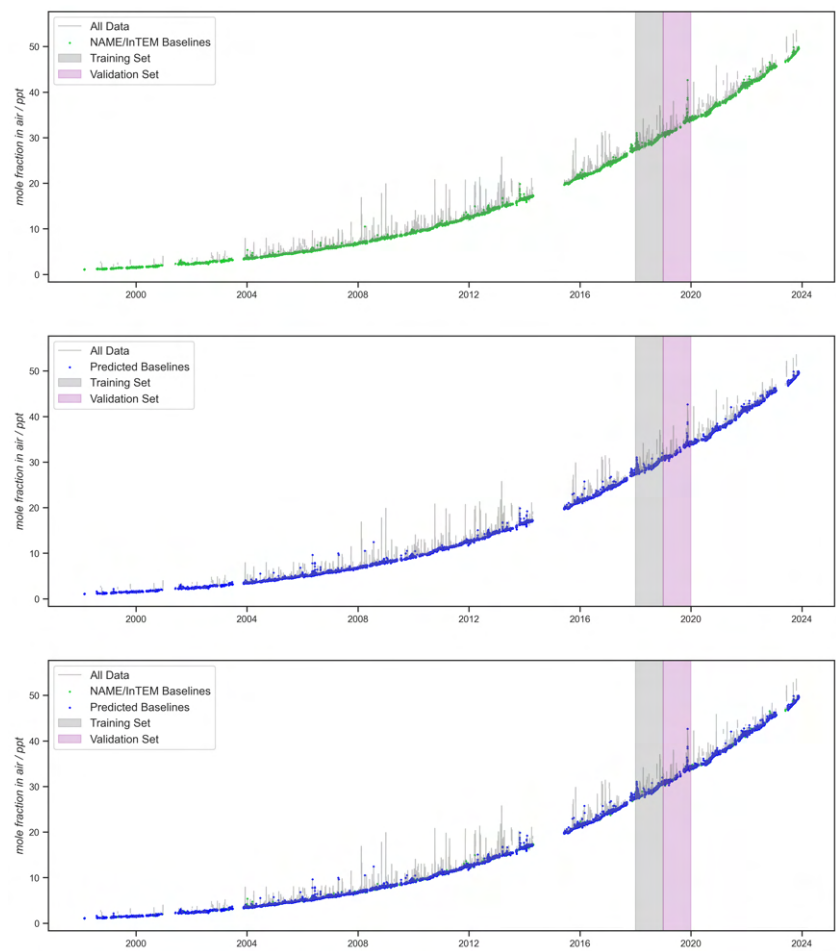
Mole fraction time series



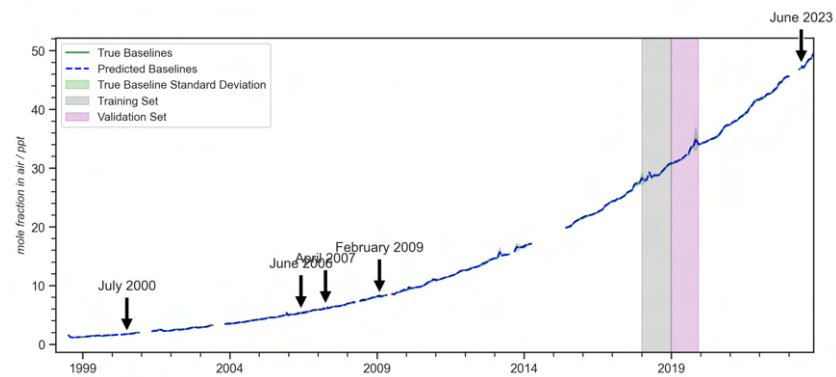
Monthly means



Mole fraction time series

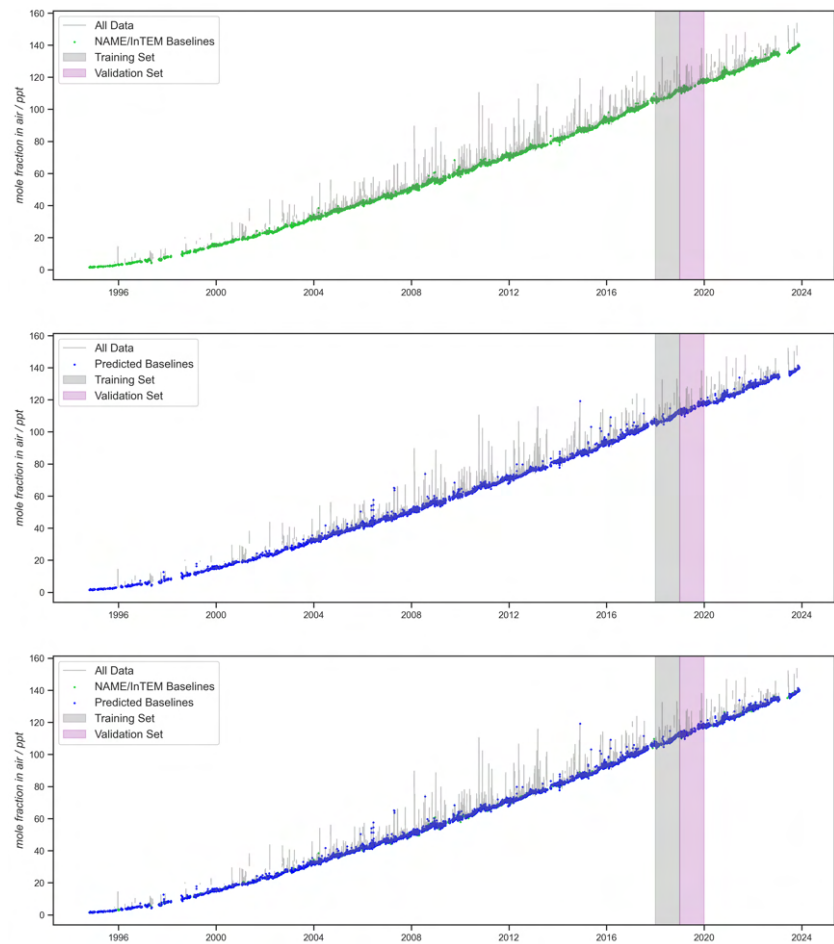


Monthly means

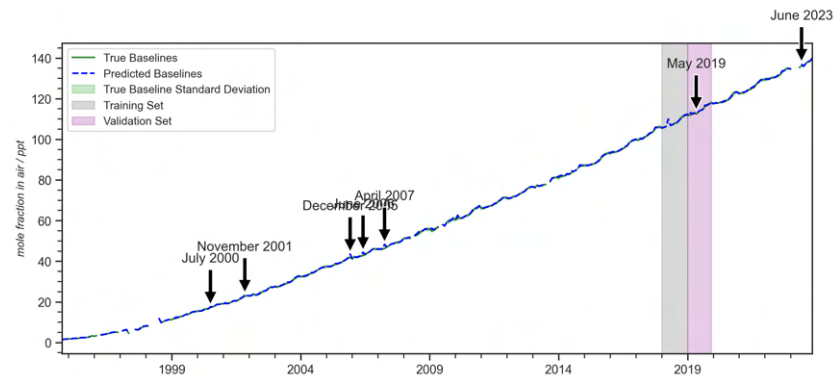


4.8.8 HFC-134a

Mole fraction time series

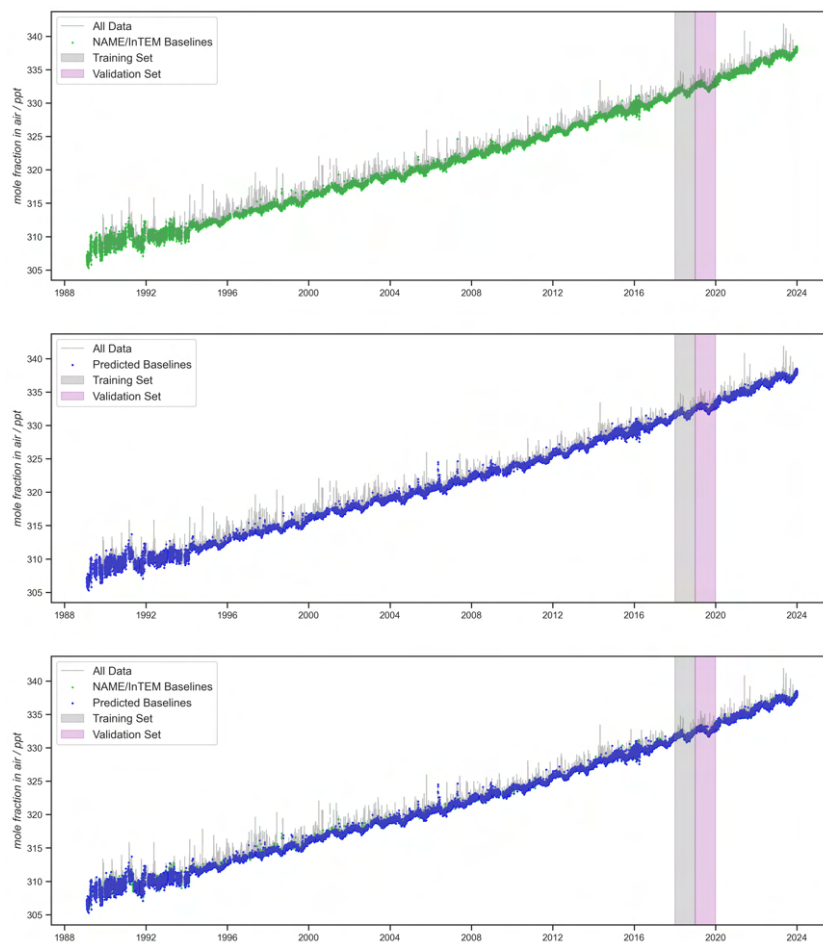


140 Monthly means

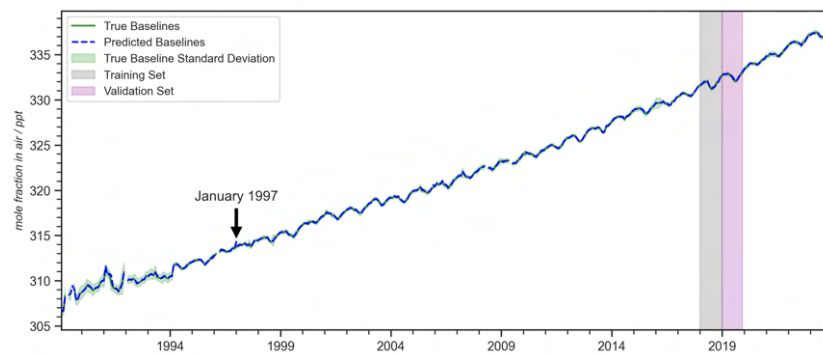


4.8.9 N₂O

Mole fraction time series

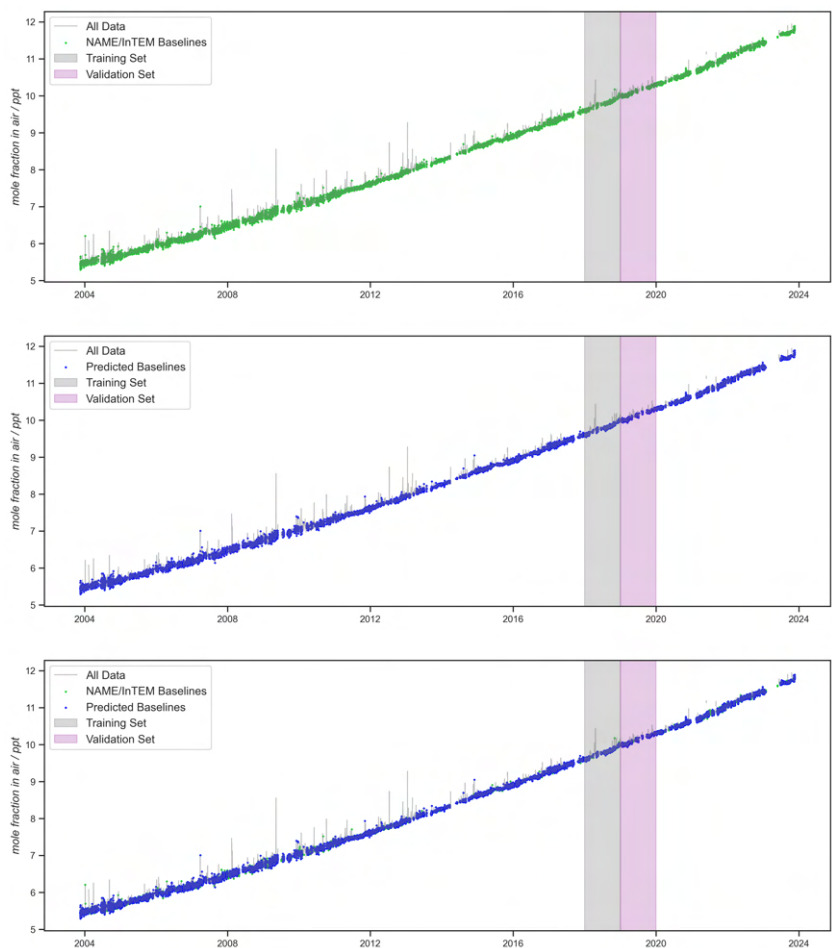


Monthly means

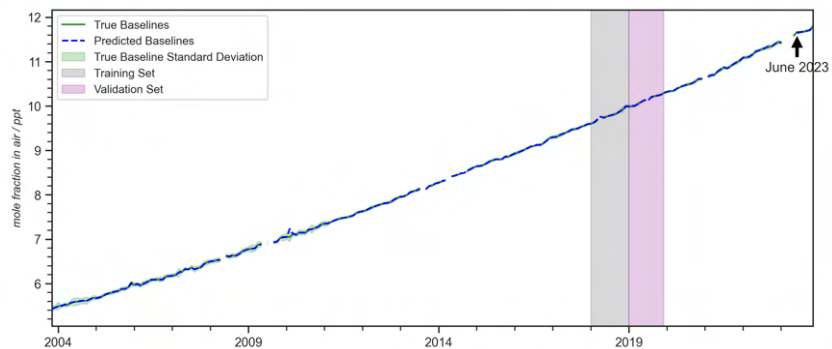


4.8.10 SF₆

145 Mole fraction time series



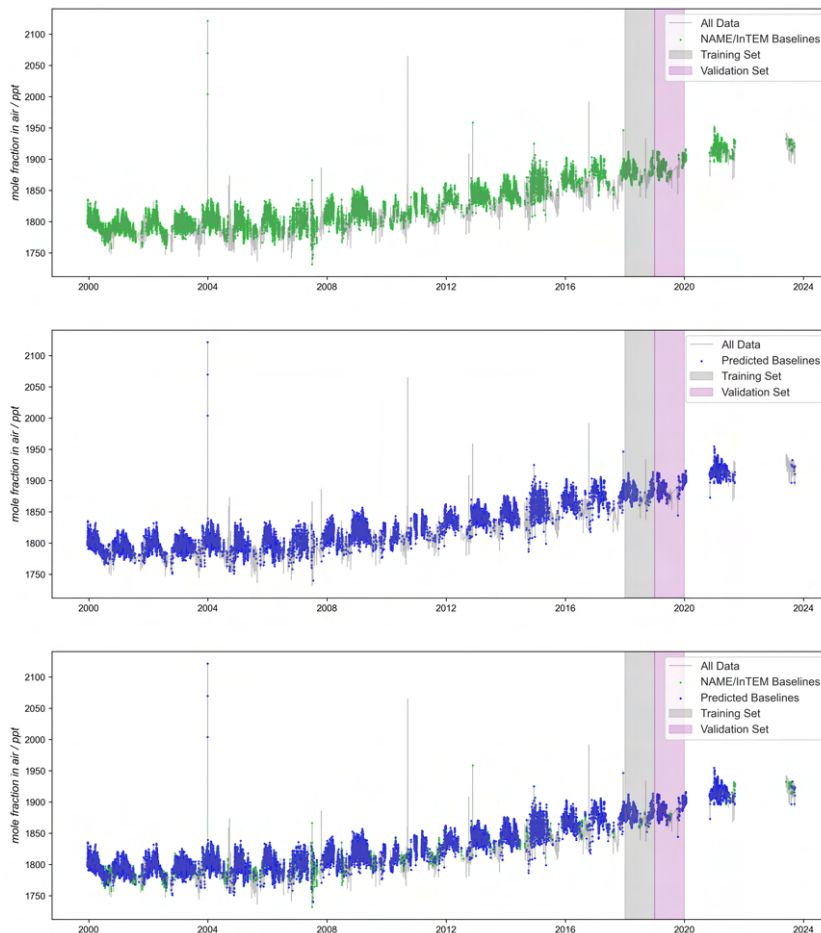
Monthly means

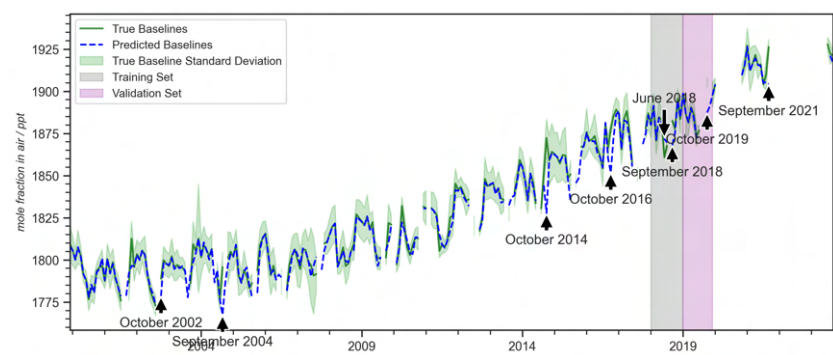


4.9 Ragged Point, Barbados

4.9.1 CH₄

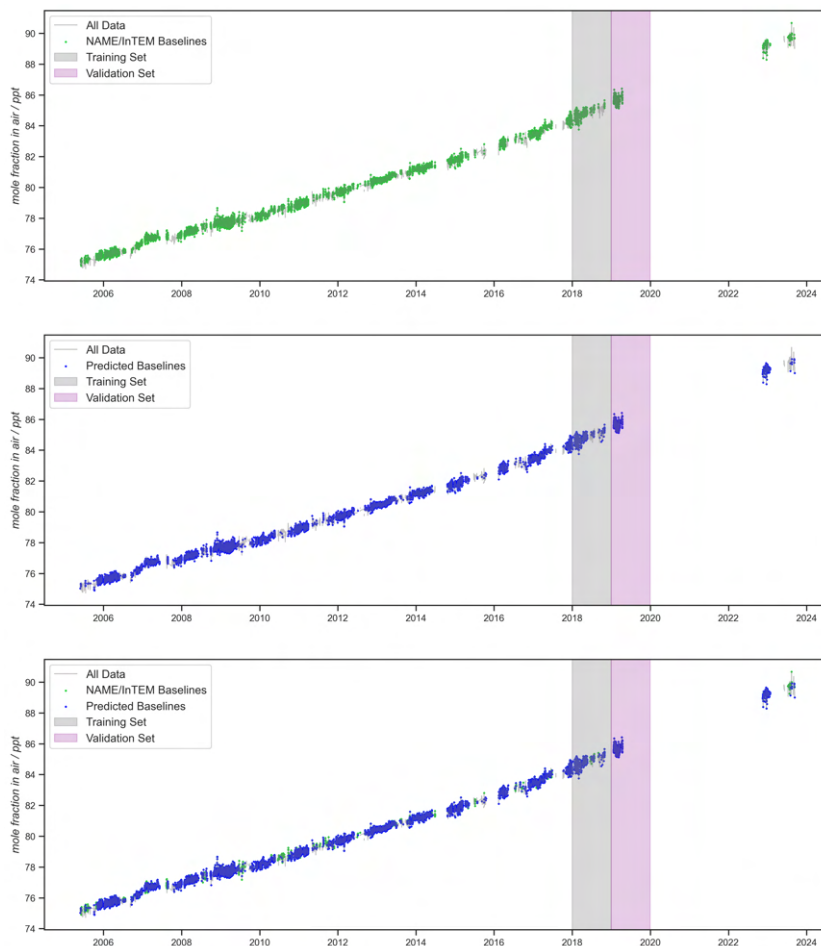
Mole fraction time series



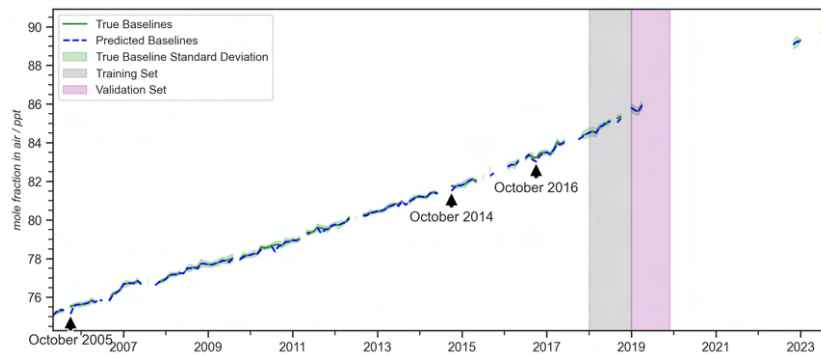


4.9.2 CF₄

Mole fraction time series

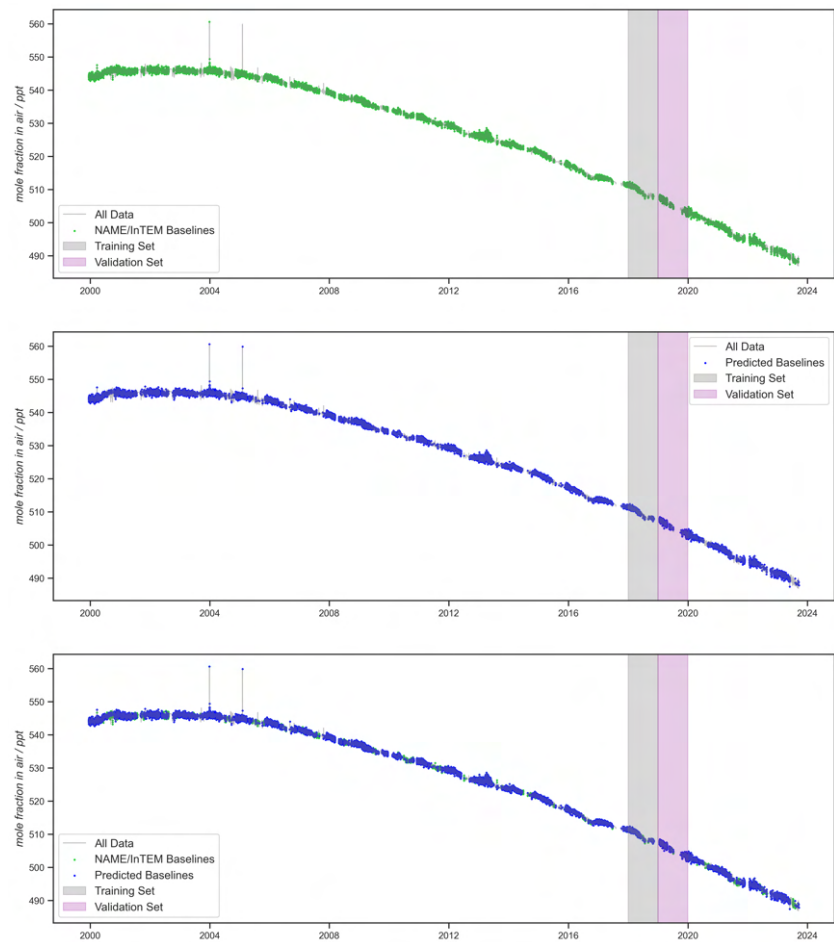


Monthly means

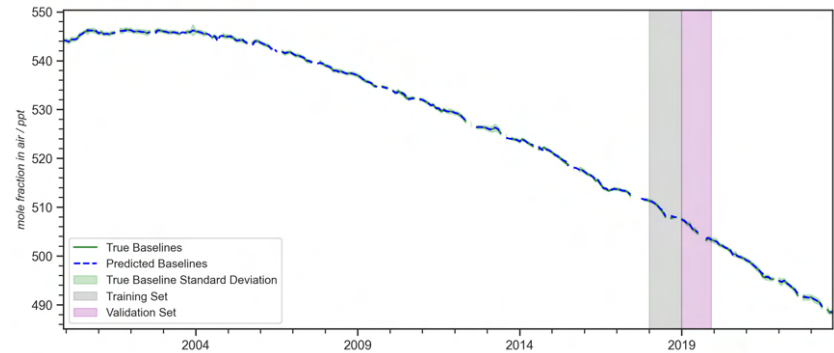


4.9.3 CFC-12

155 Mole fraction time series

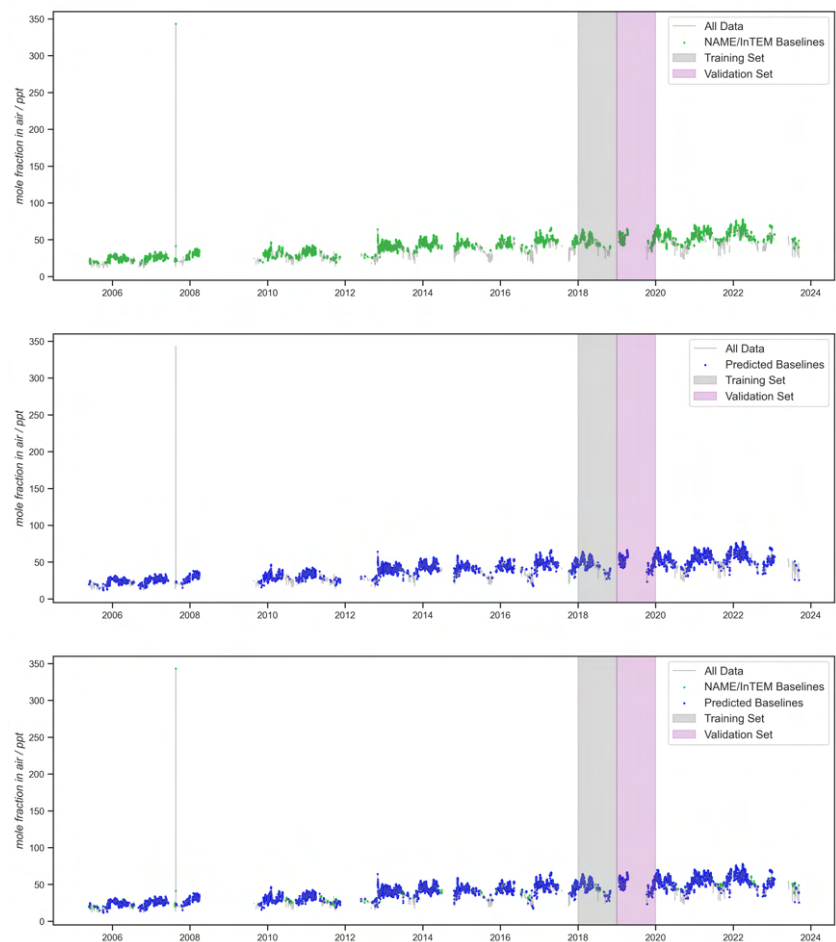


Monthly means

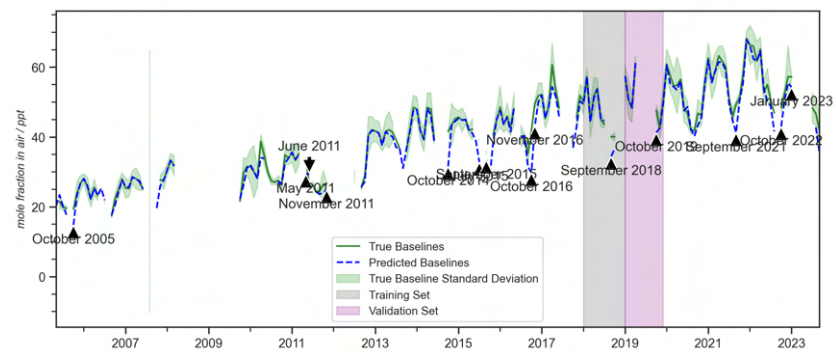


4.9.4 CH₂Cl₂

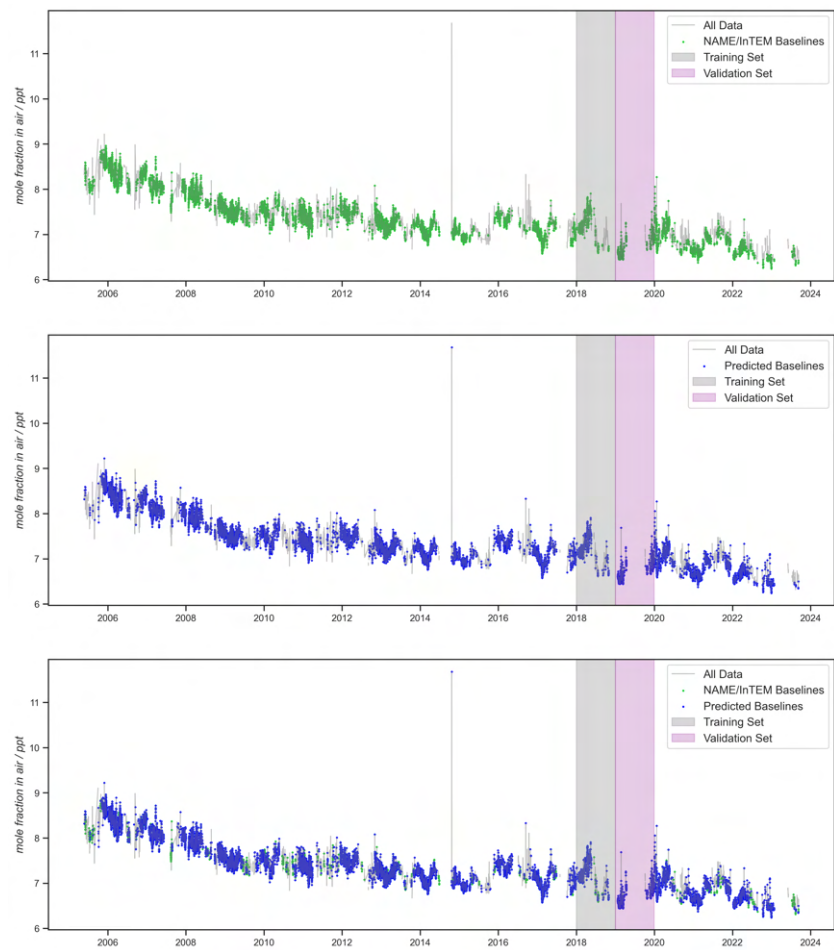
Mole fraction time series



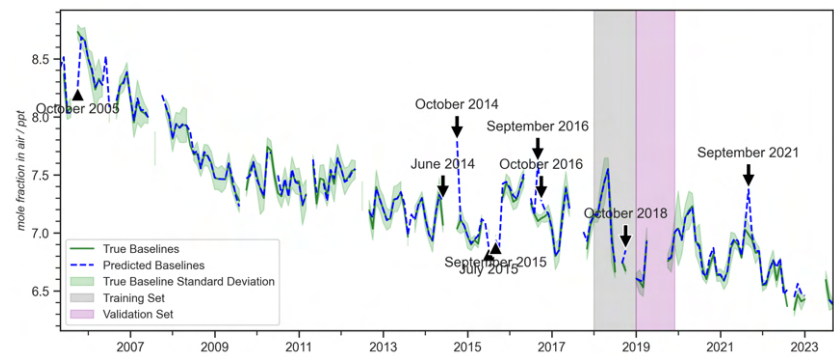
Monthly means



Mole fraction time series

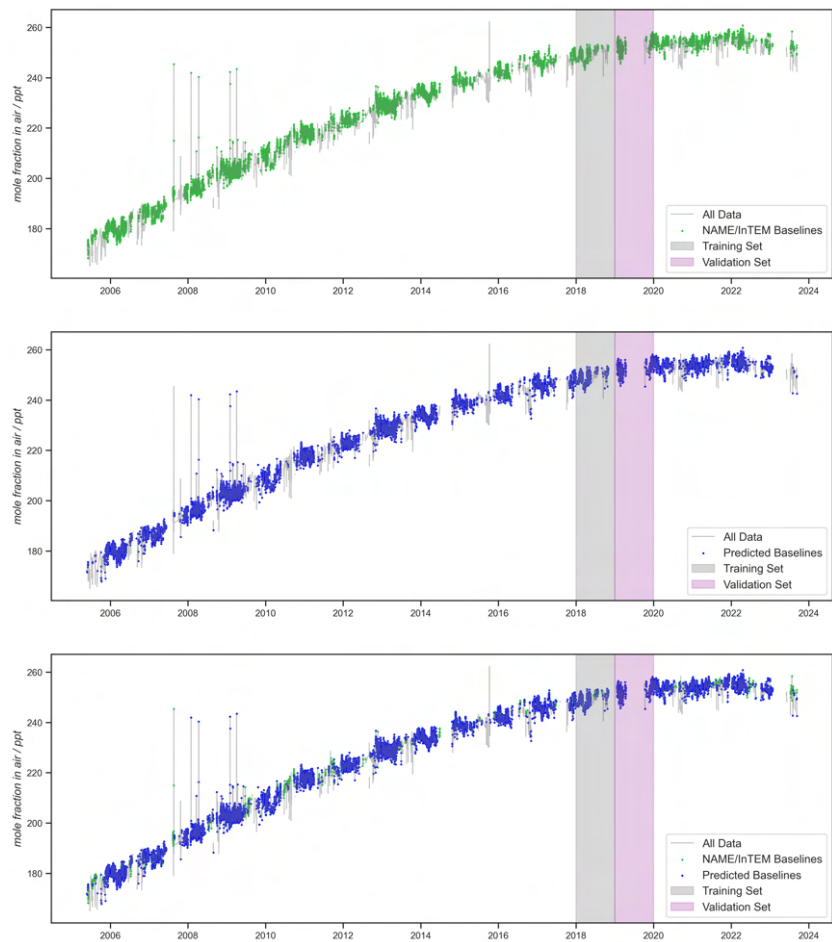


Monthly means

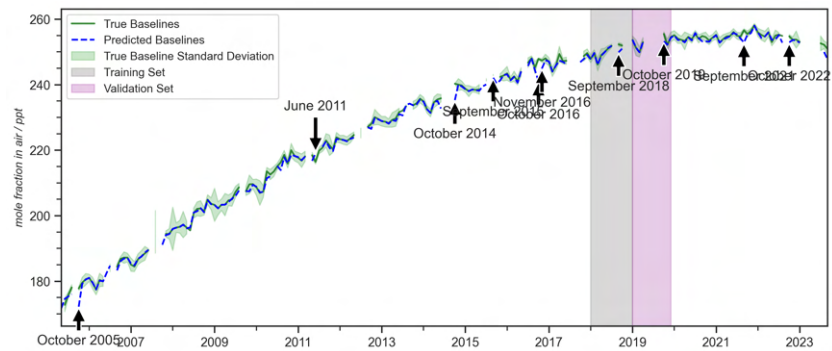


4.9.6 HCFC-22

Mole fraction time series

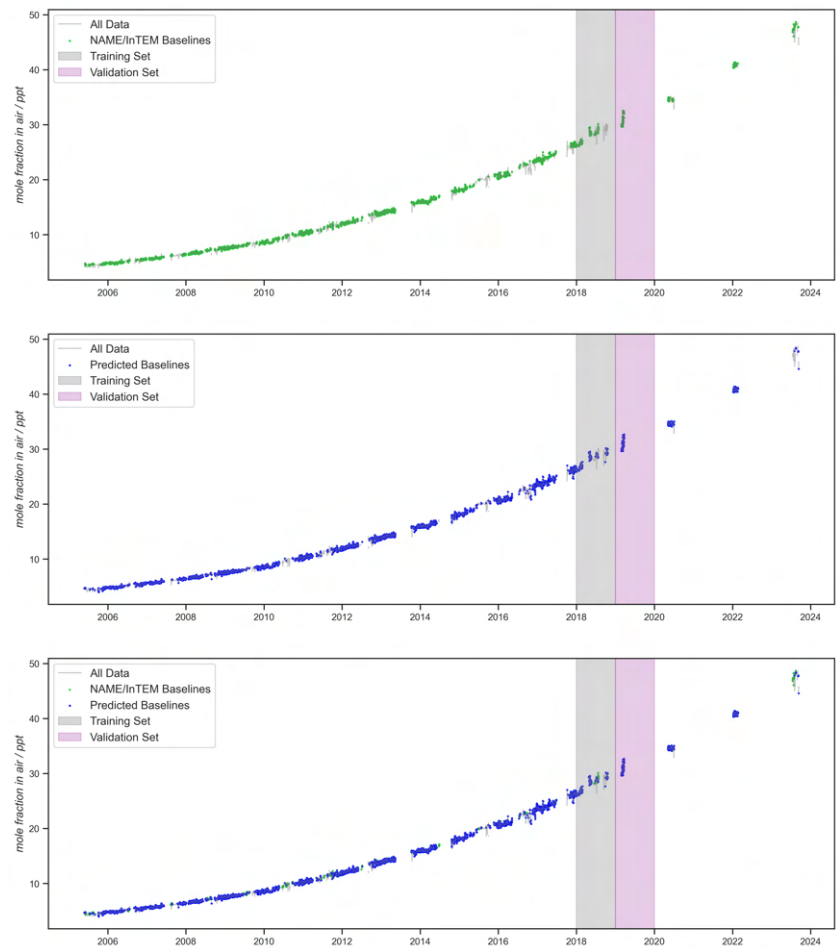


165 Monthly means

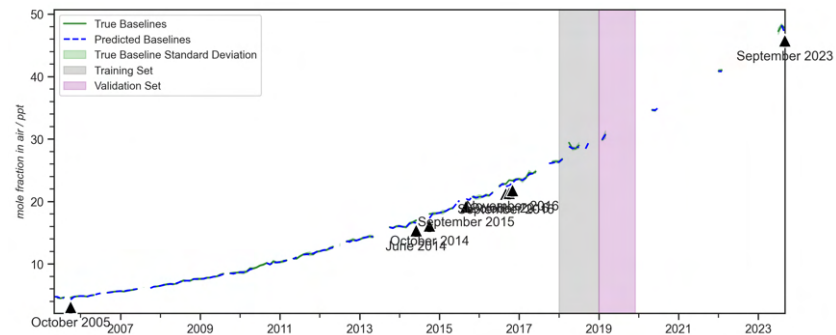


4.9.7 HFC-125

Mole fraction time series

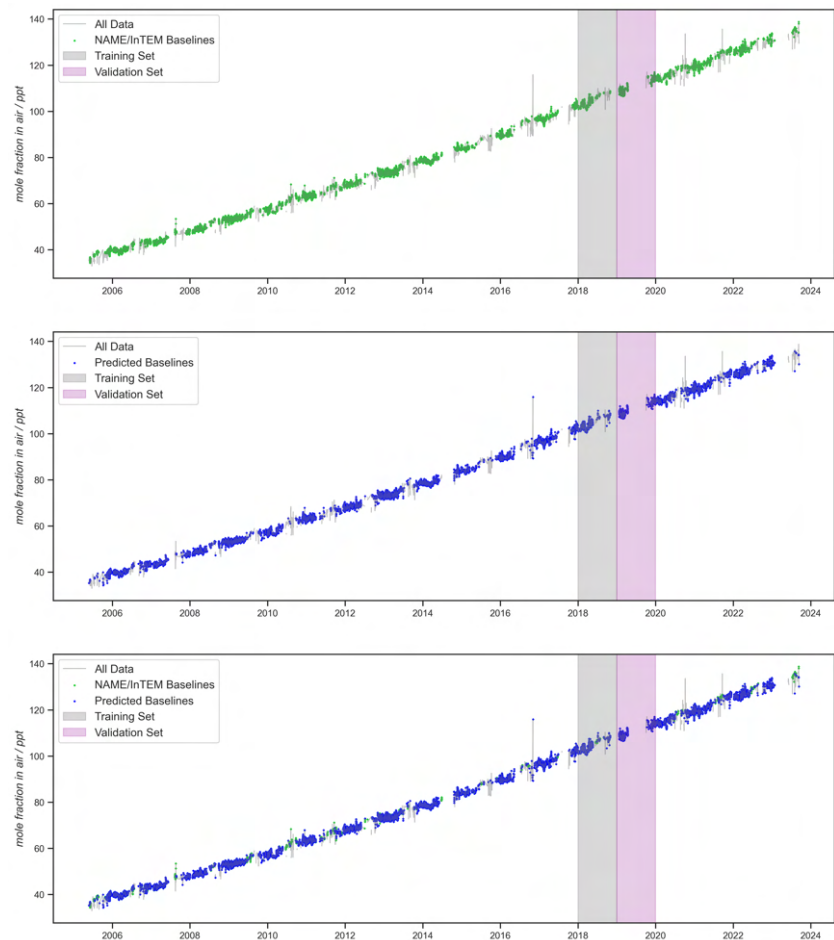


Monthly means

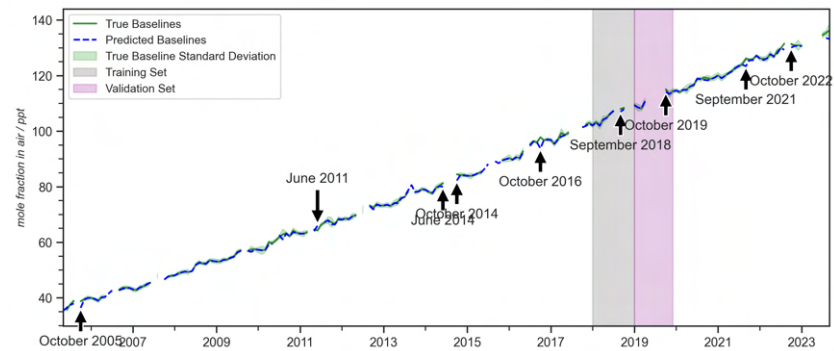


4.9.8 HFC-134a

170 Mole fraction time series

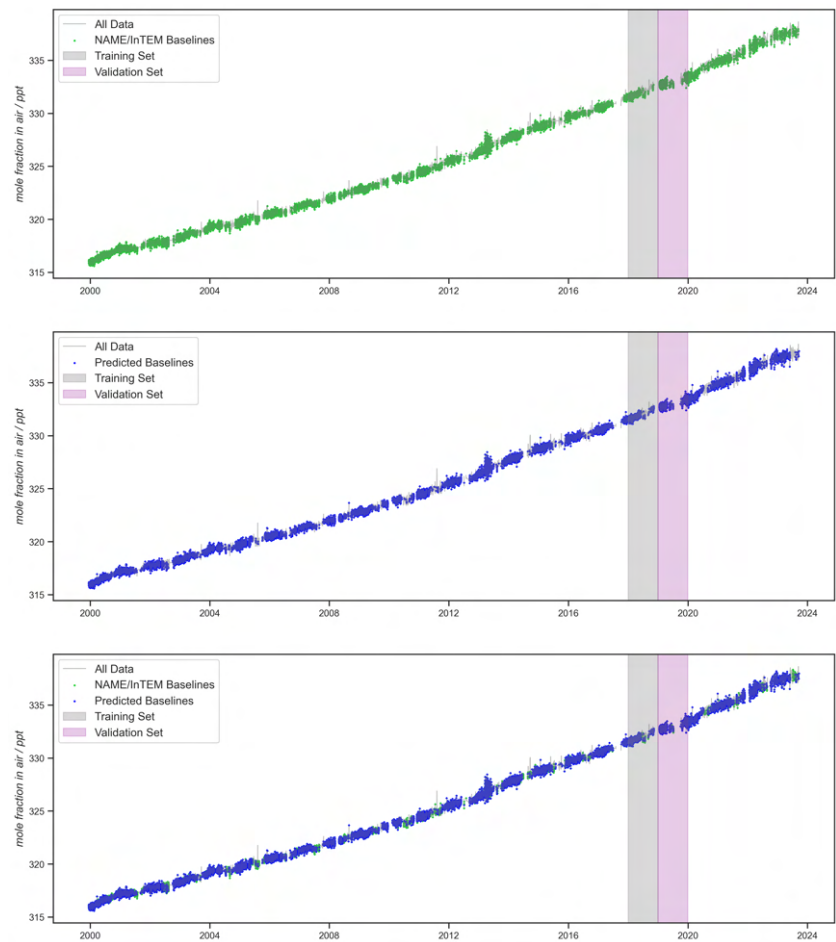


Monthly means

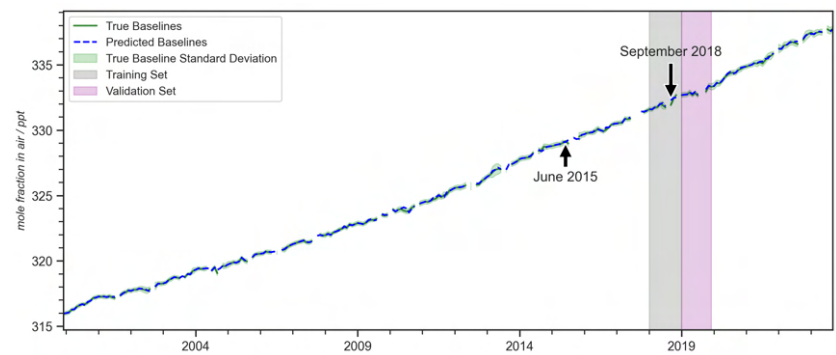


4.9.9 N₂O

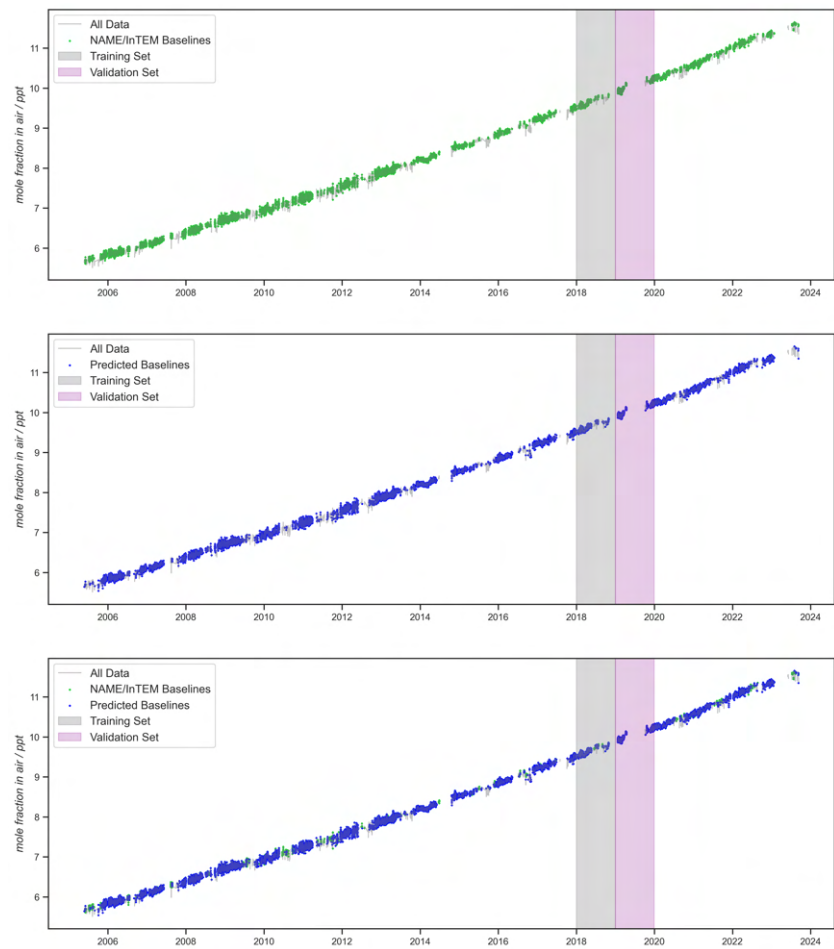
Mole fraction time series



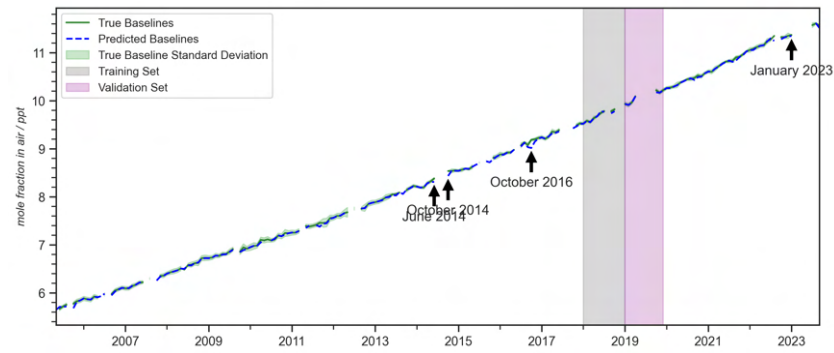
Monthly means



Mole fraction time series



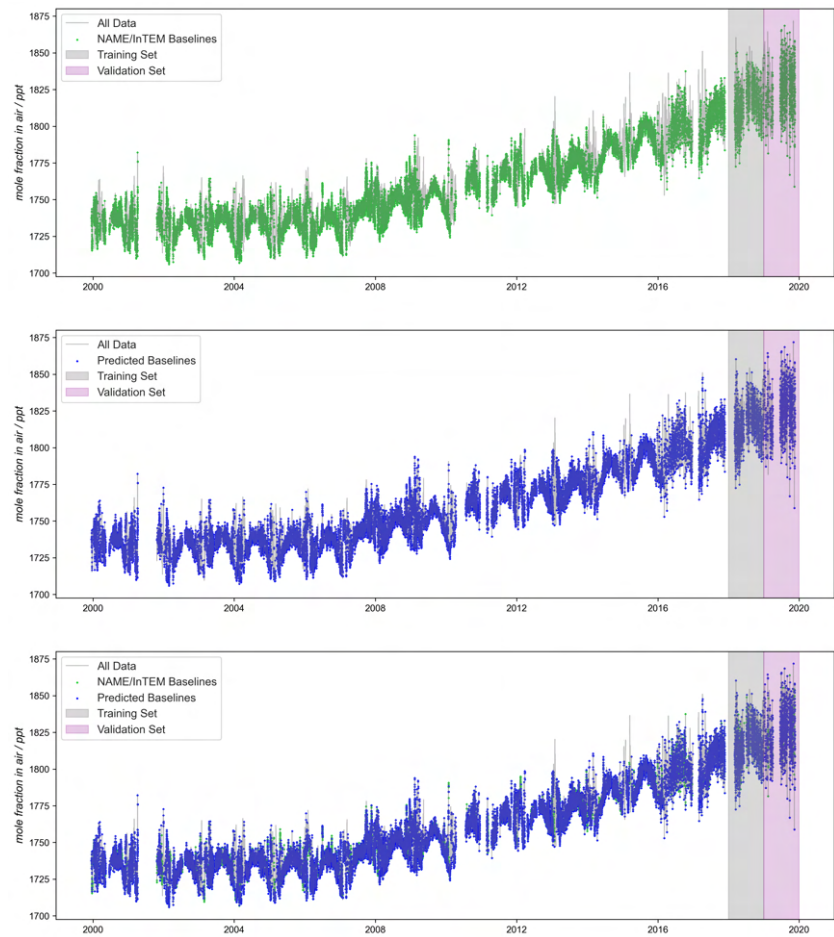
Monthly means



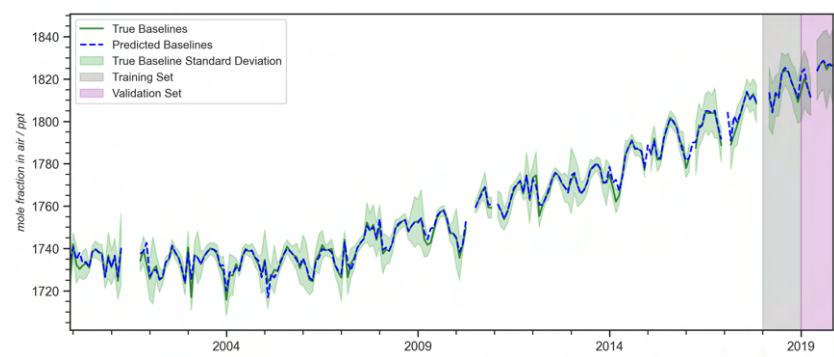
4.10 Cape Matatula, American Samoa

4.10.1 CH₄

180 Mole fraction time series

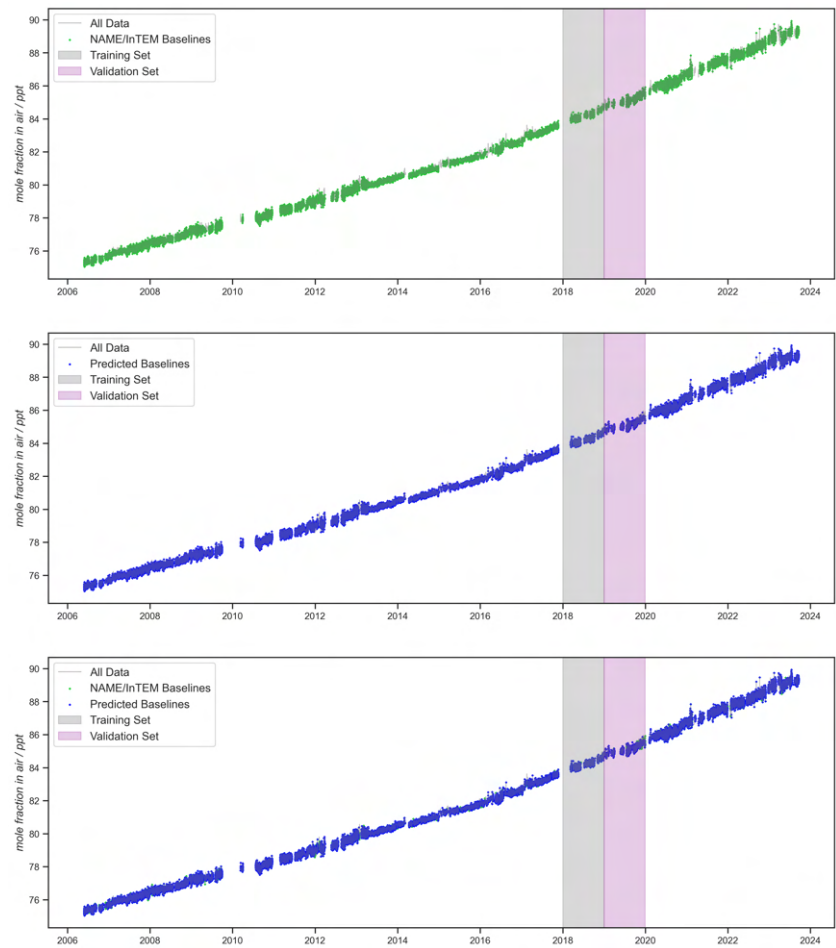


Monthly means

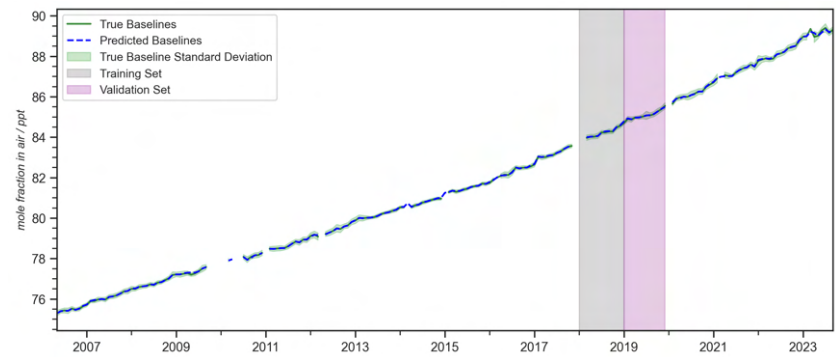


4.10.2 CF₄

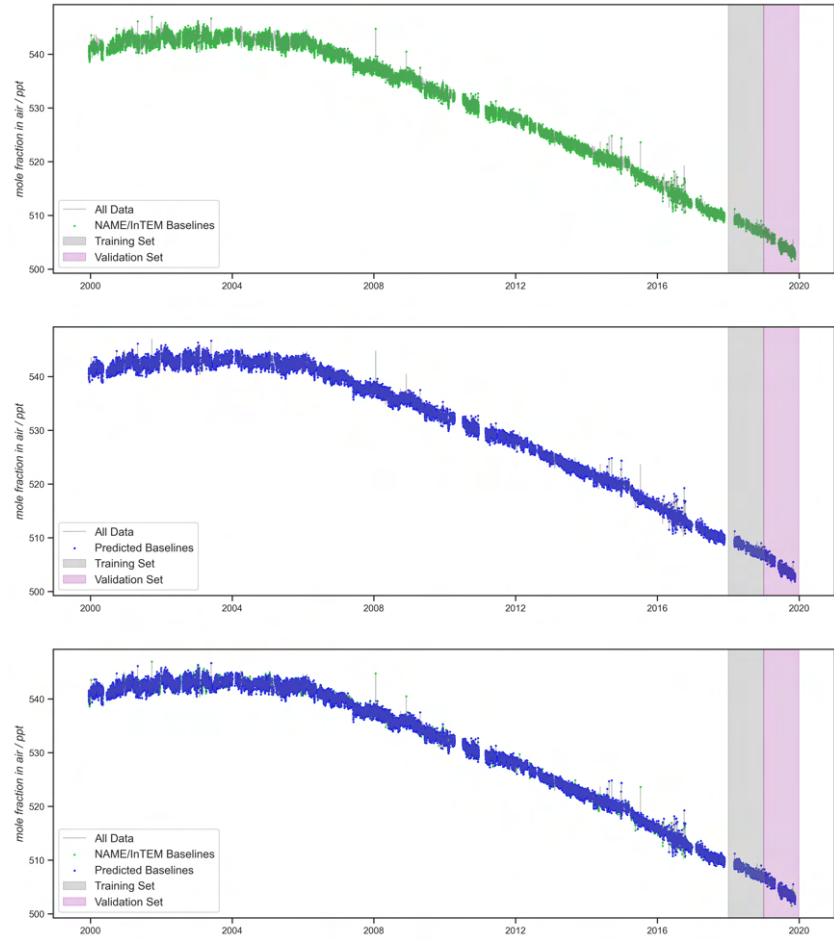
Mole fraction time series



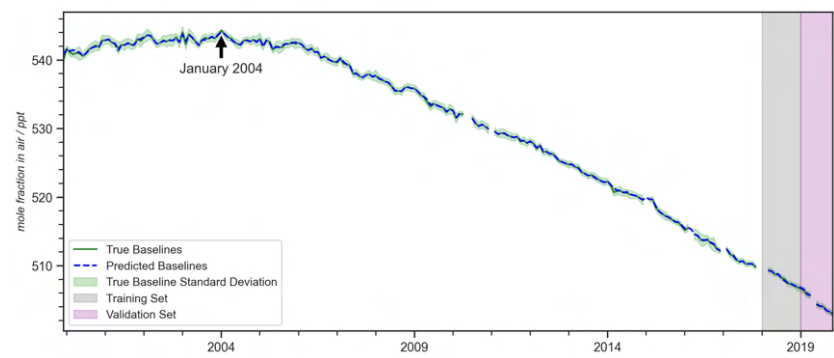
Monthly means



Mole fraction time series

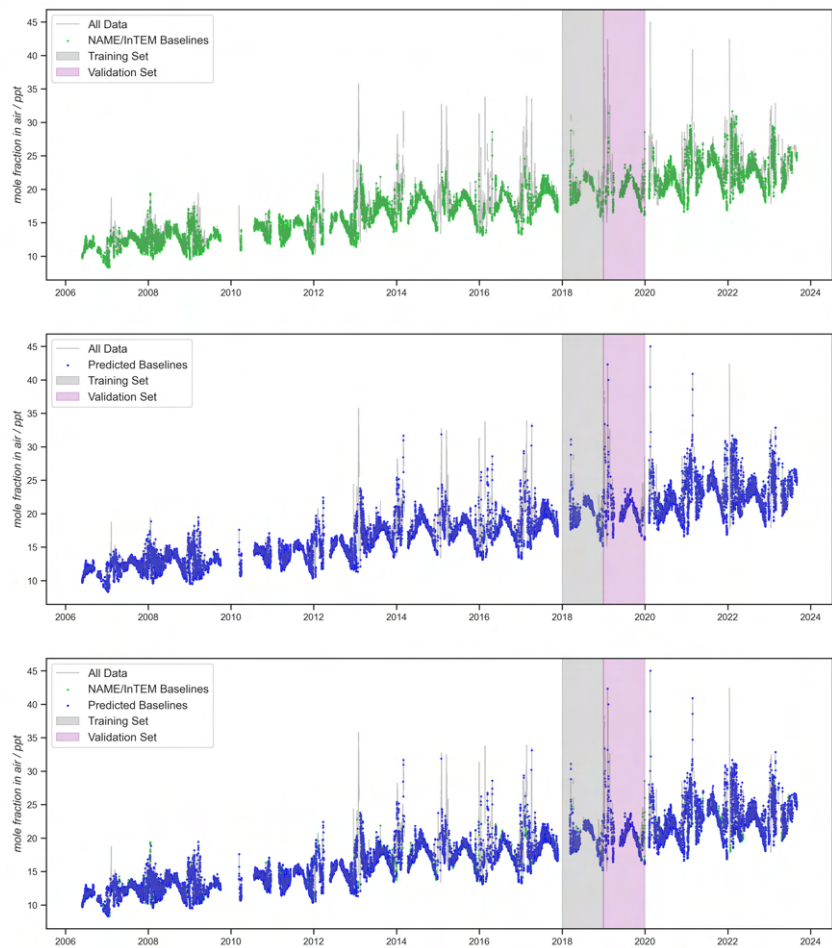


Monthly means

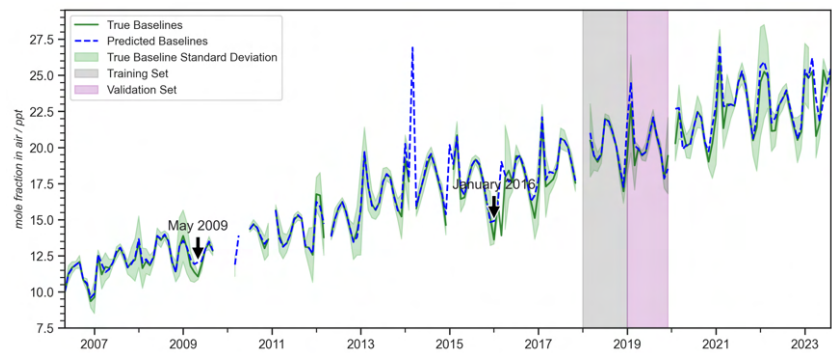


4.10.4 CH₂Cl₂

Mole fraction time series

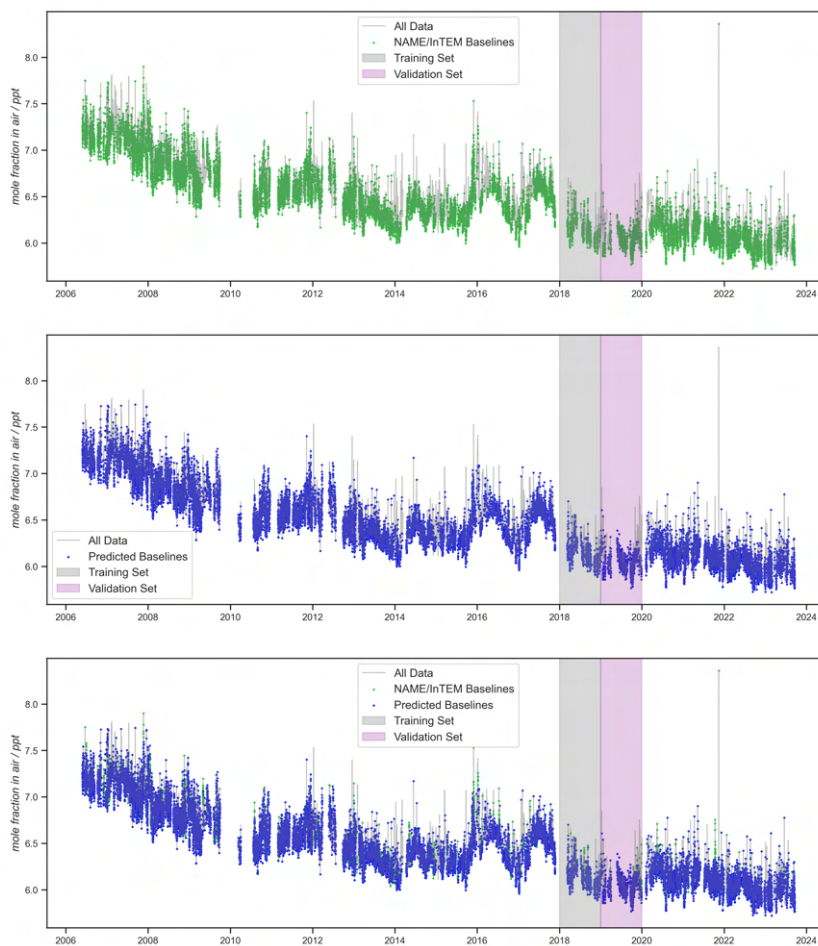


190 Monthly means

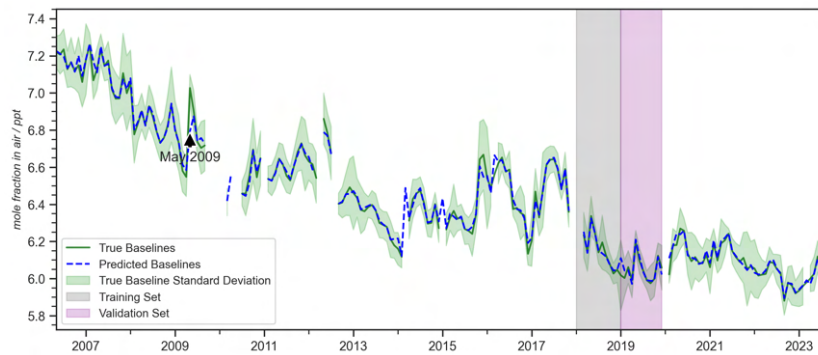


4.10.5 CH₃Br

Mole fraction time series

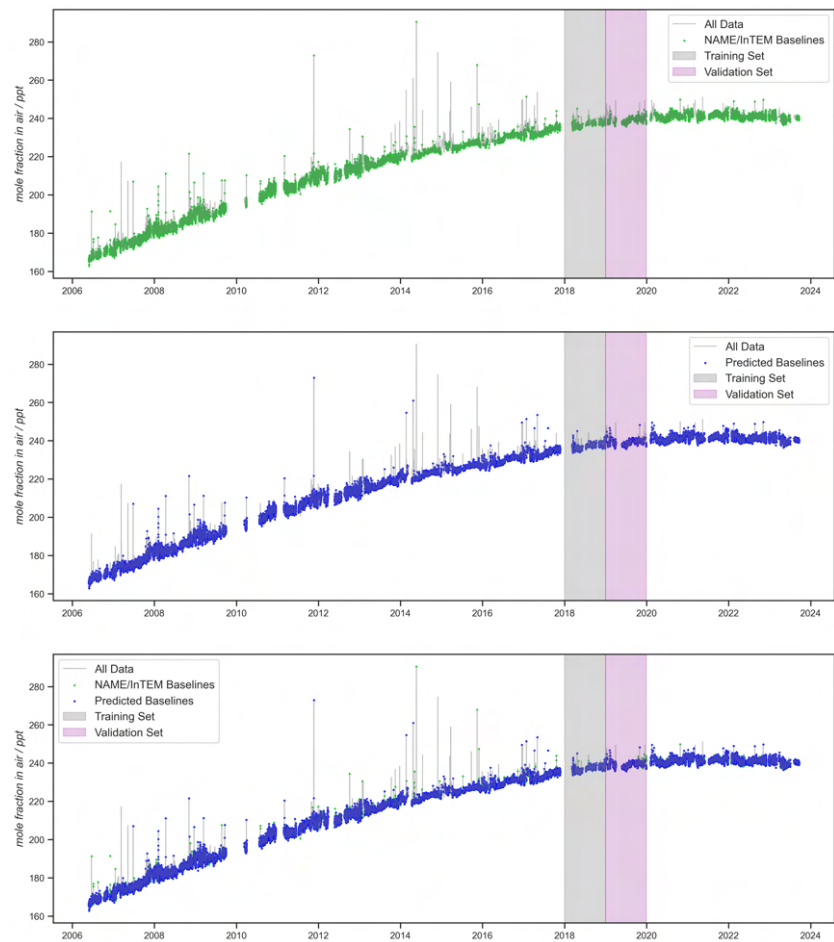


Monthly means

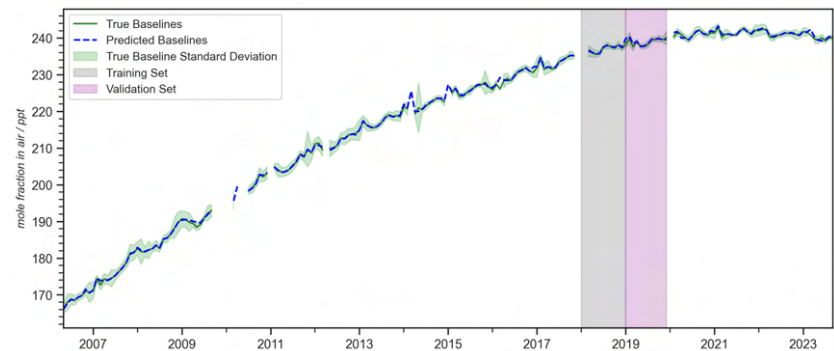


4.10.6 HCFC-22

195 Mole fraction time series

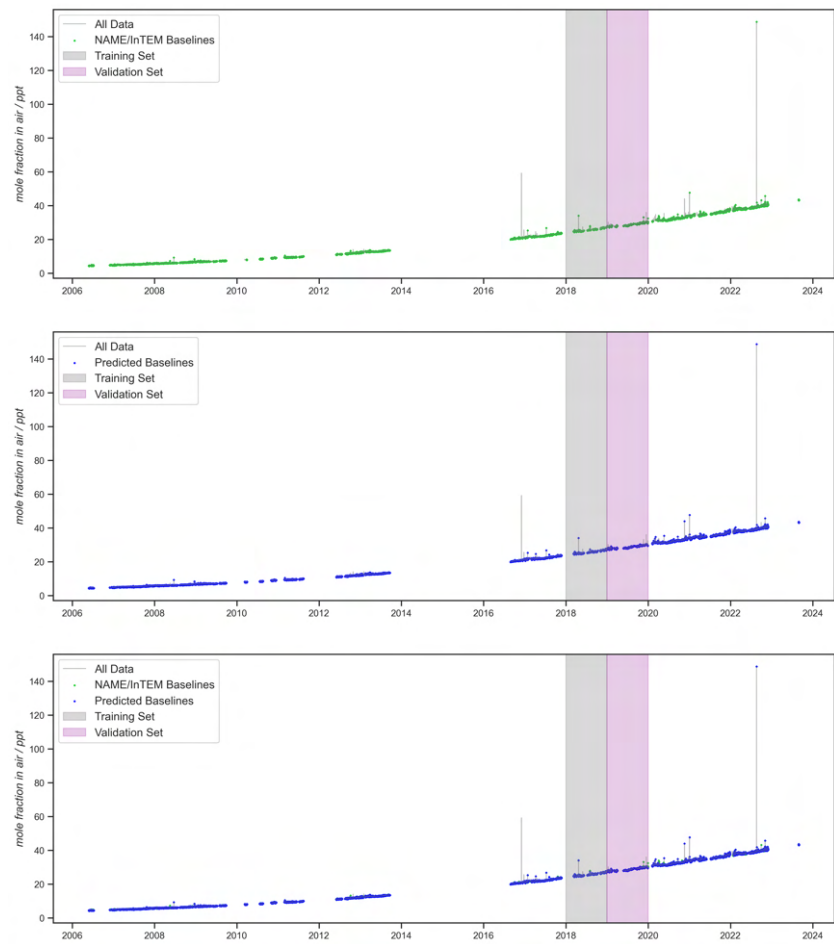


Monthly means

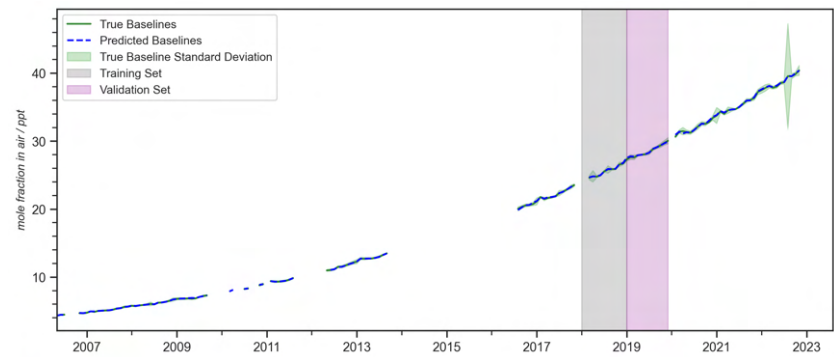


4.10.7 HFC-125

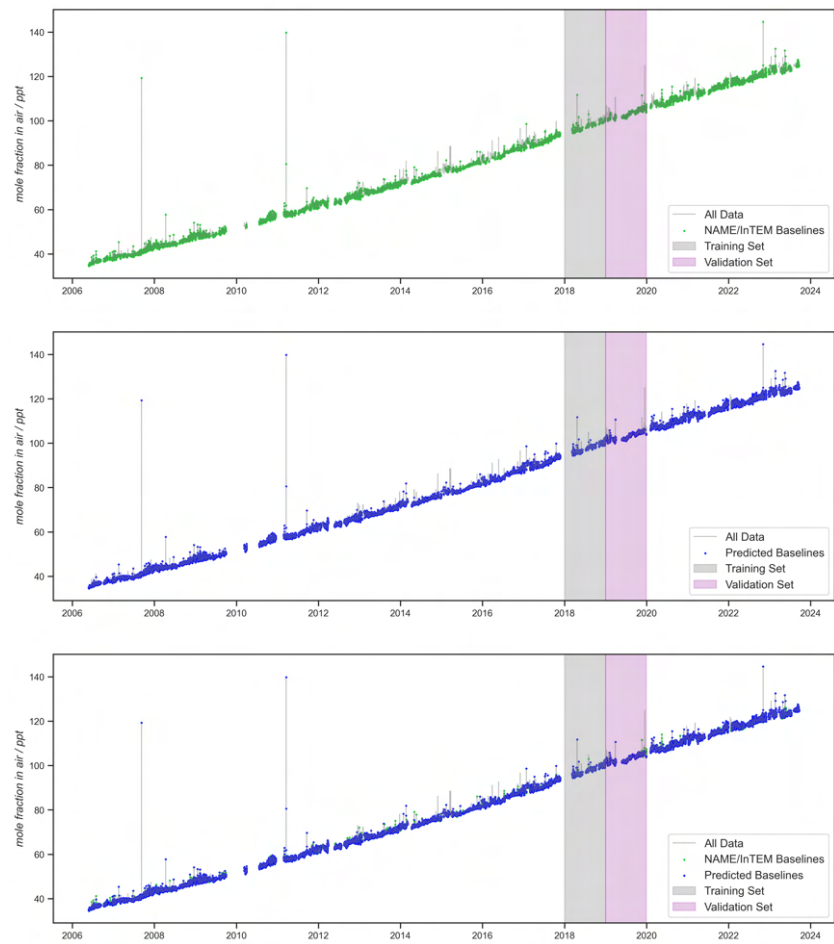
Mole fraction time series



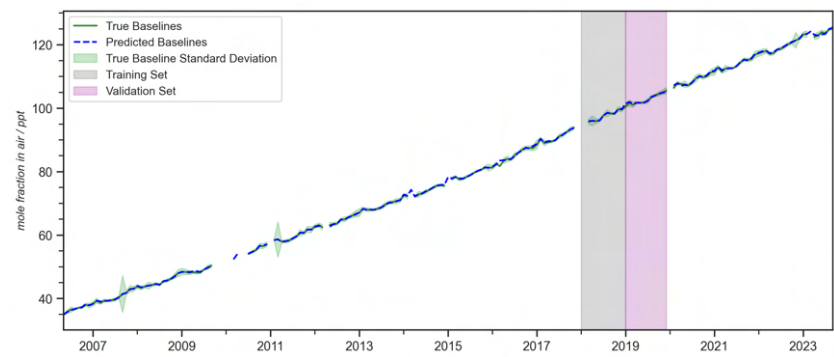
Monthly means



Mole fraction time series

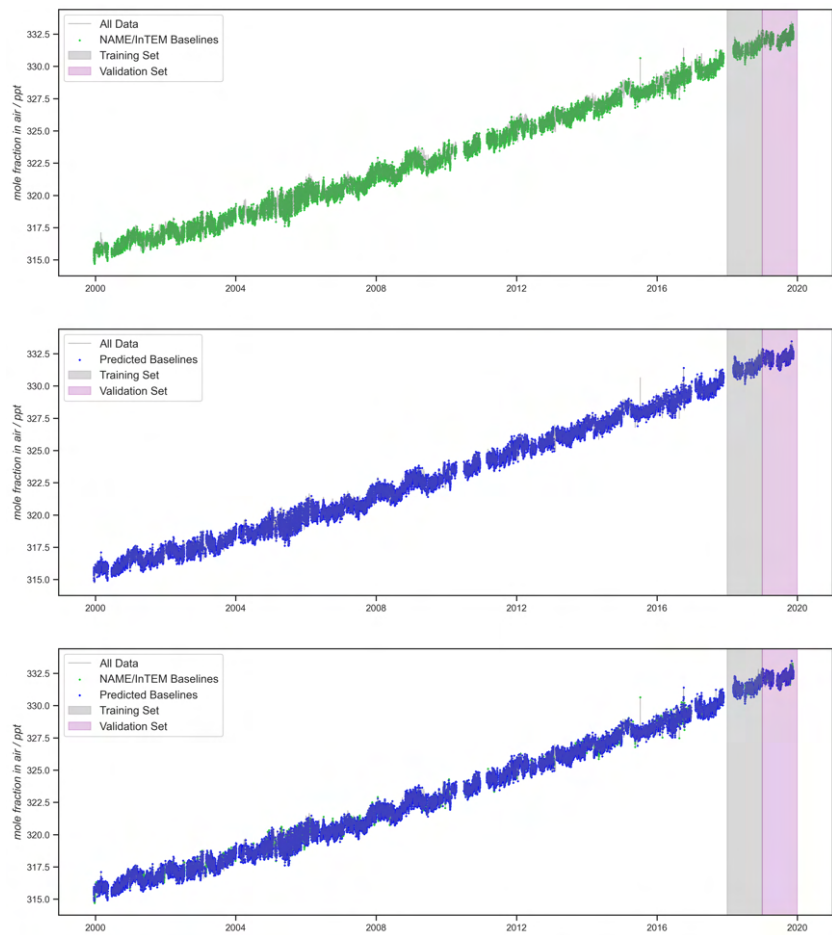


Monthly means

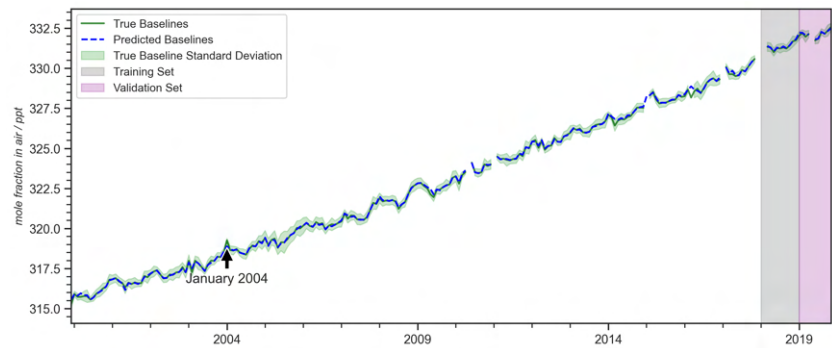


4.10.9 N₂O

Mole fraction time series

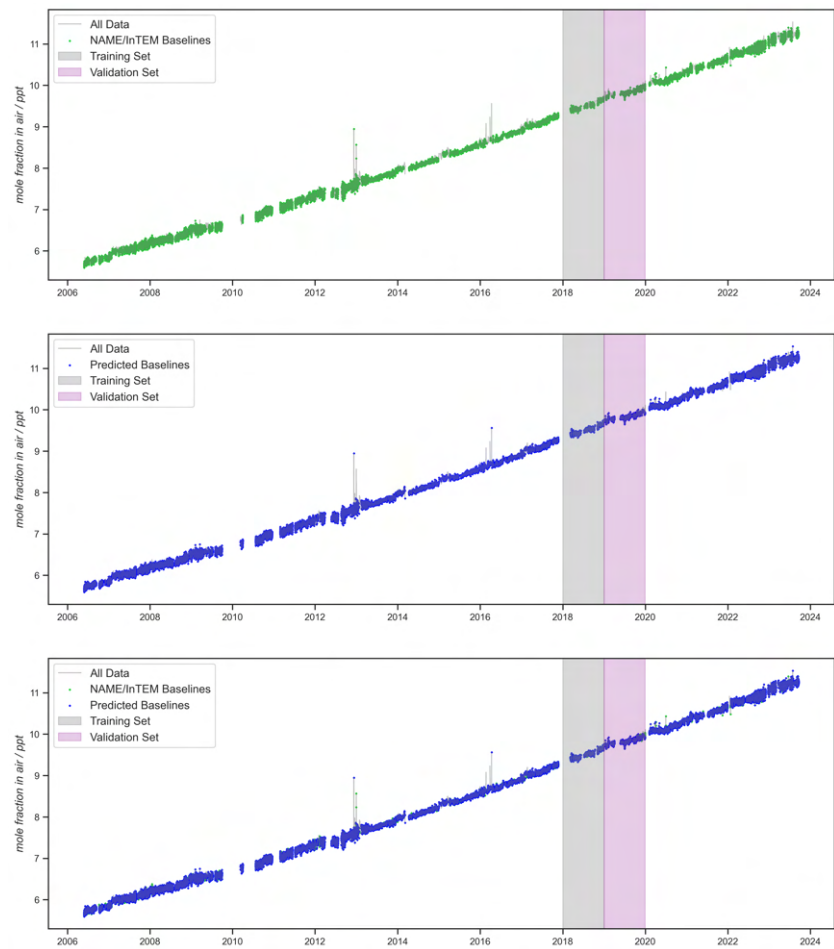


205 Monthly means

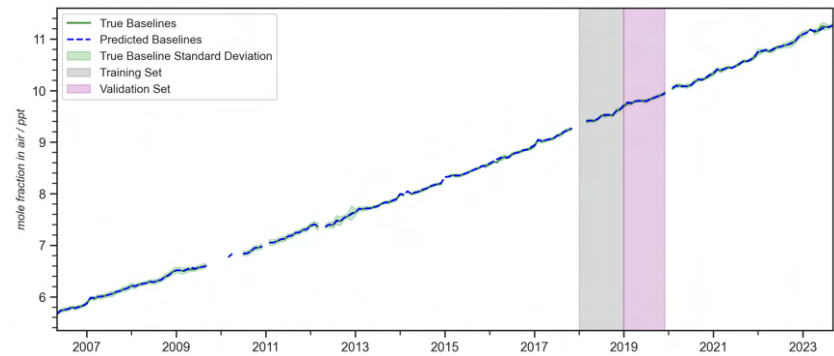


4.10.10 SF₆

Mole fraction time series



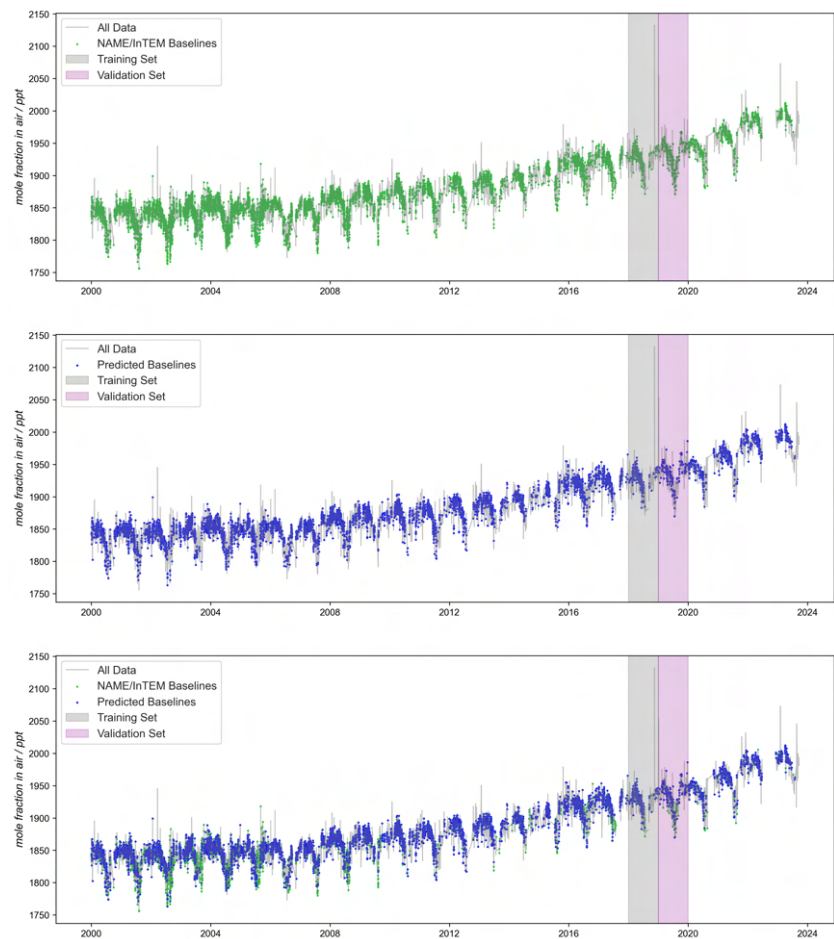
Monthly means



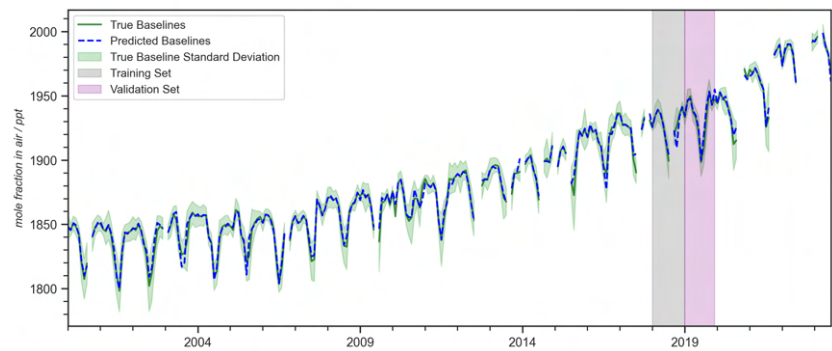
4.11 Trinidad Head, USA

210 4.11.1 CH₄

Mole fraction time series

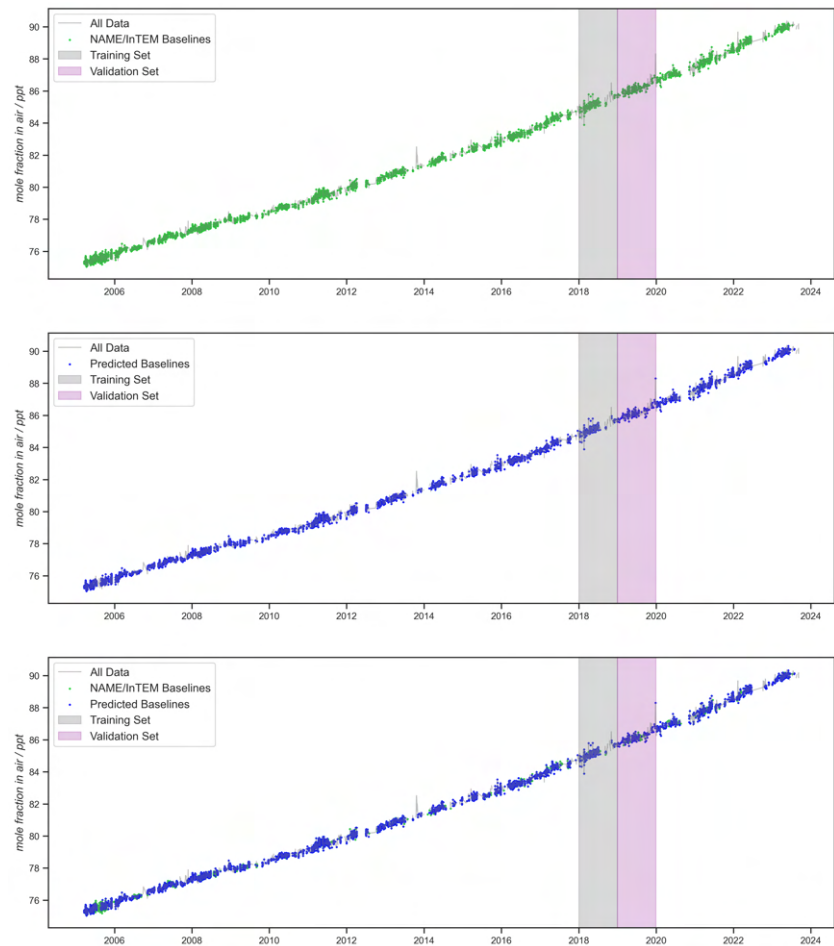


Monthly means

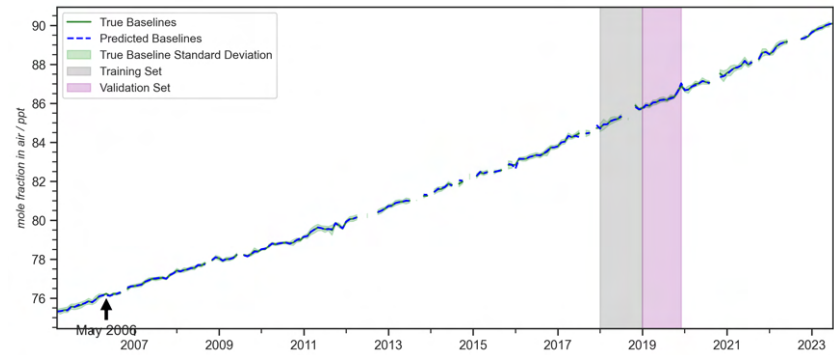


4.11.2 CF₄

Mole fraction time series

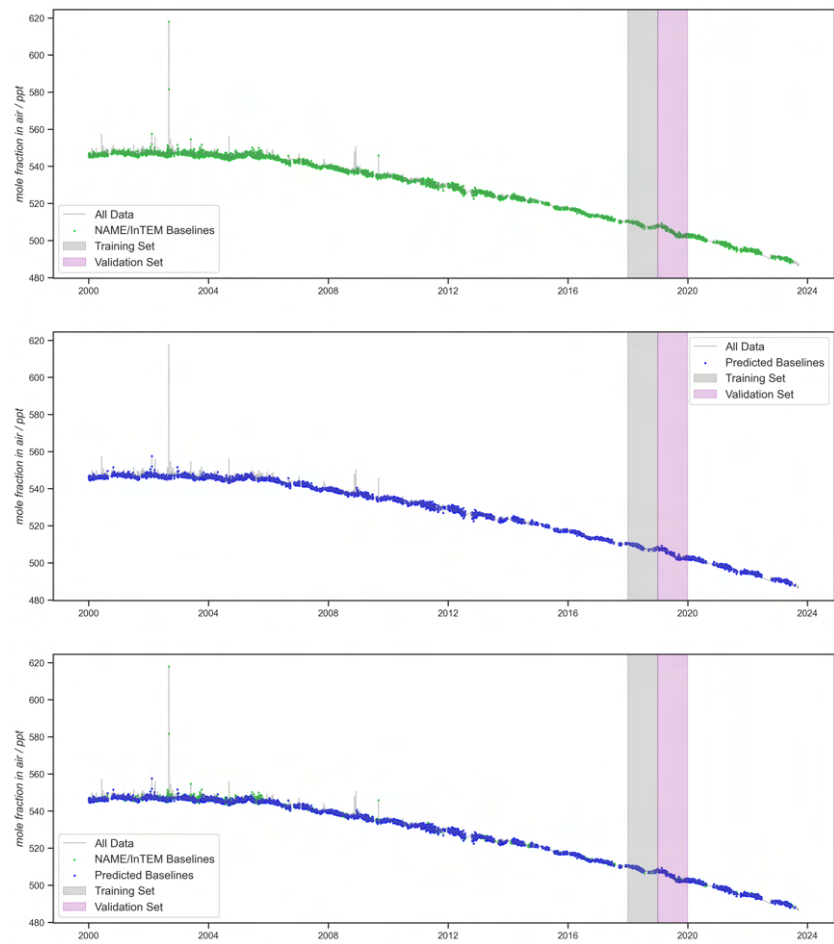


215 Monthly means

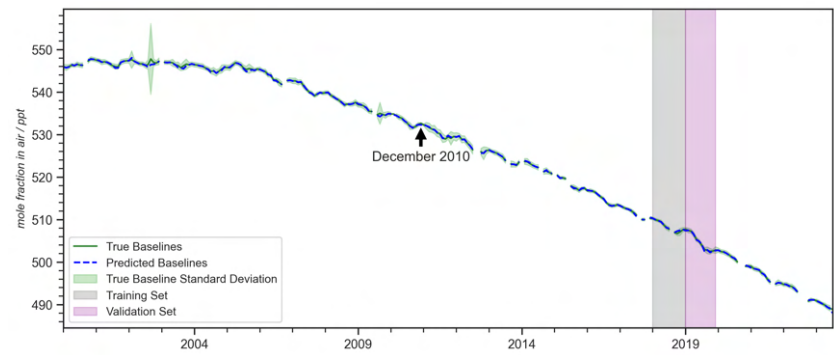


4.11.3 CFC-12

Mole fraction time series

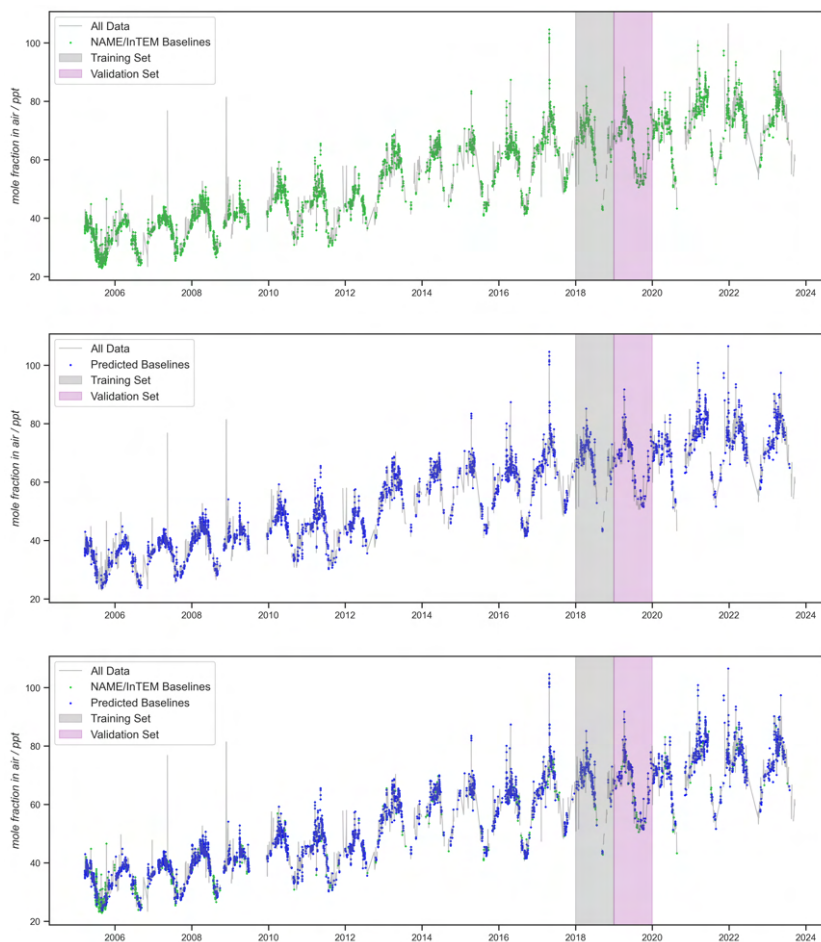


Monthly means

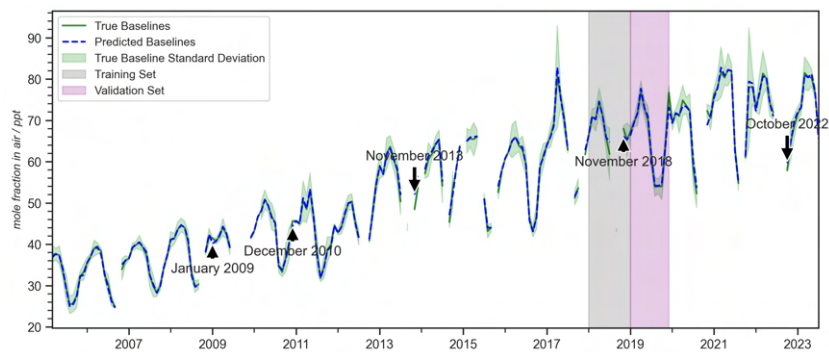


4.11.4 CH₂Cl₂

220 Mole fraction time series

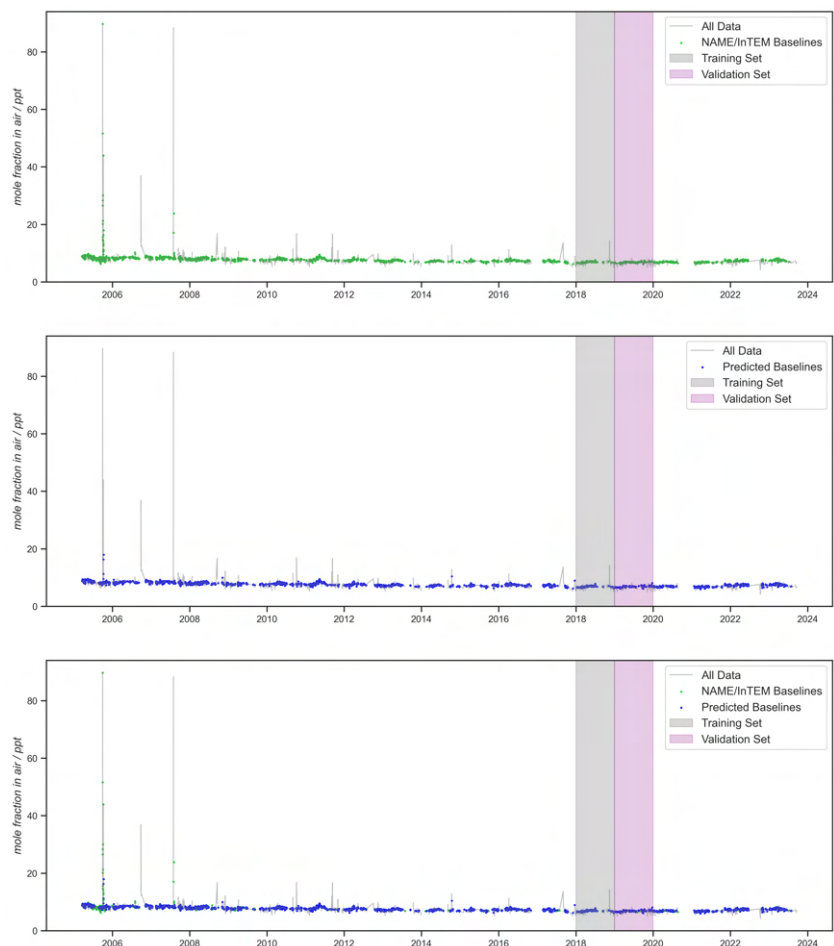


Monthly means

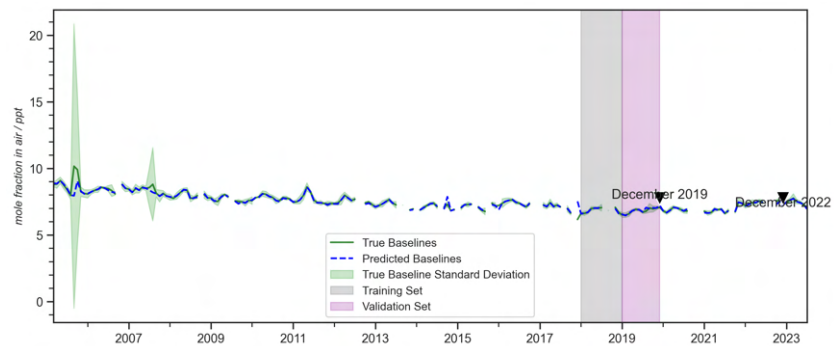


4.11.5 CH₃Br

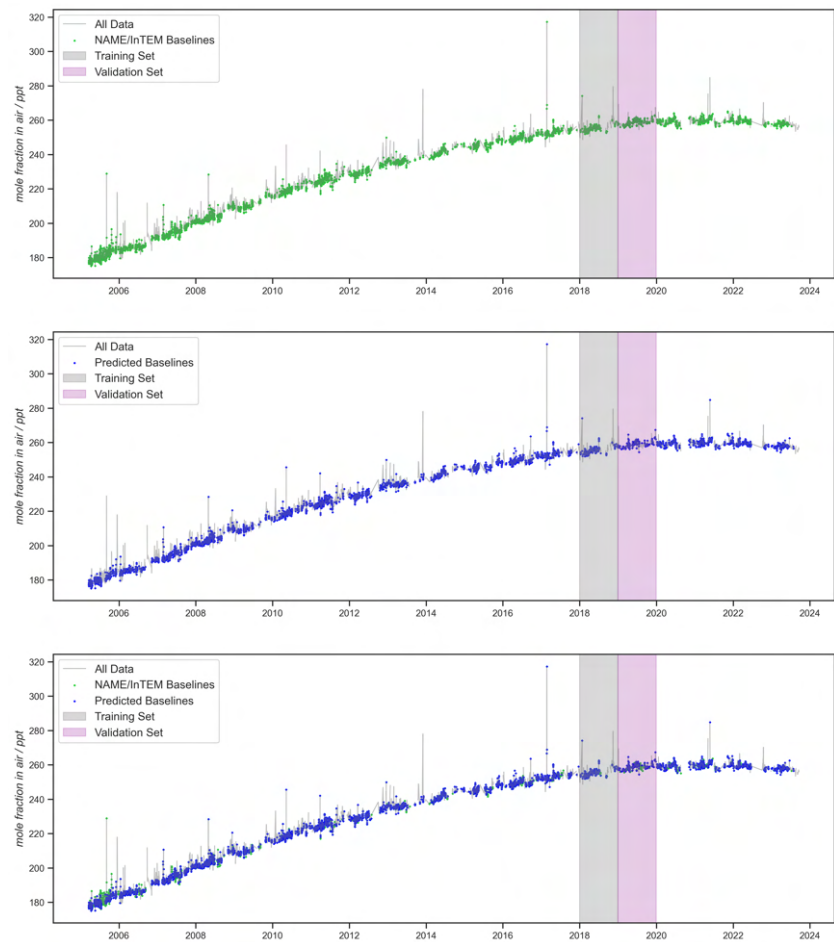
Mole fraction time series



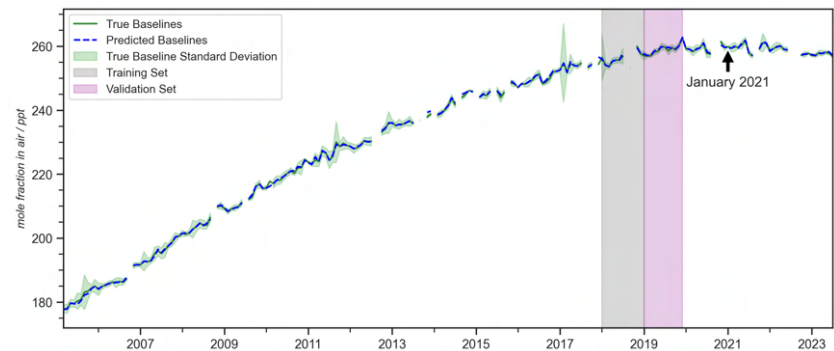
Monthly means



Mole fraction time series

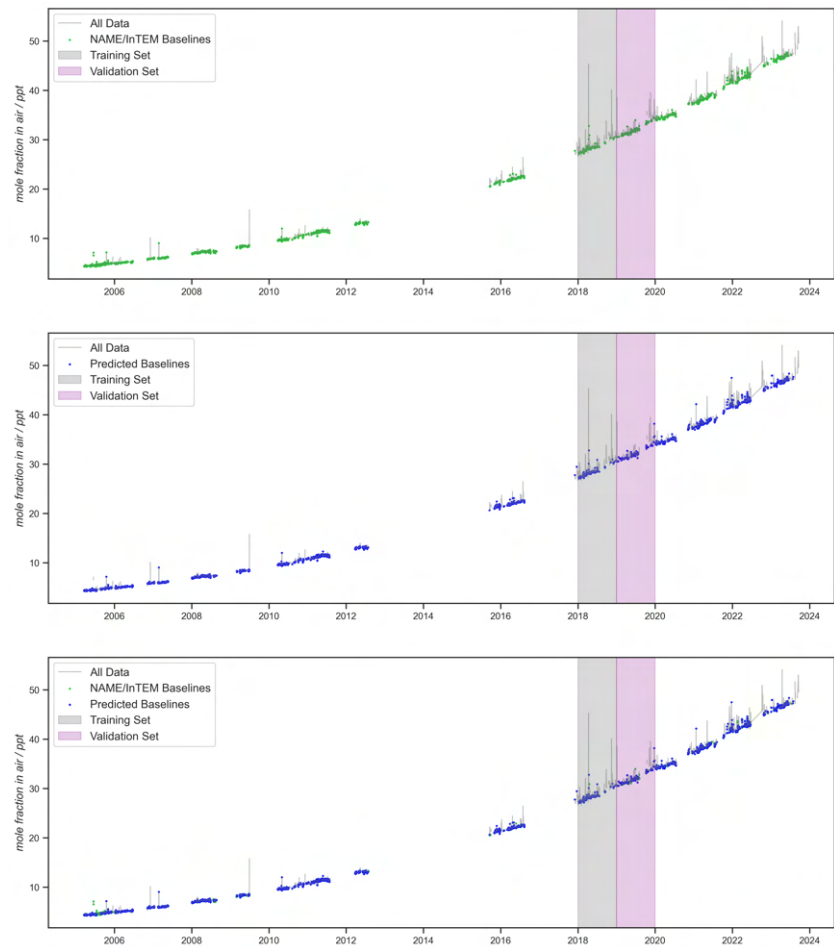


Monthly means

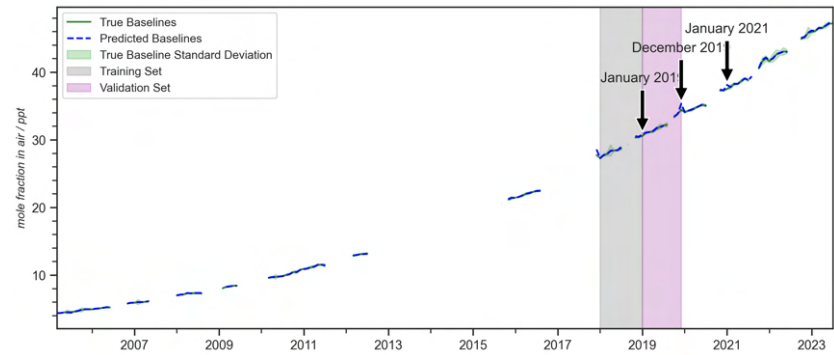


4.11.7 HFC-125

Mole fraction time series

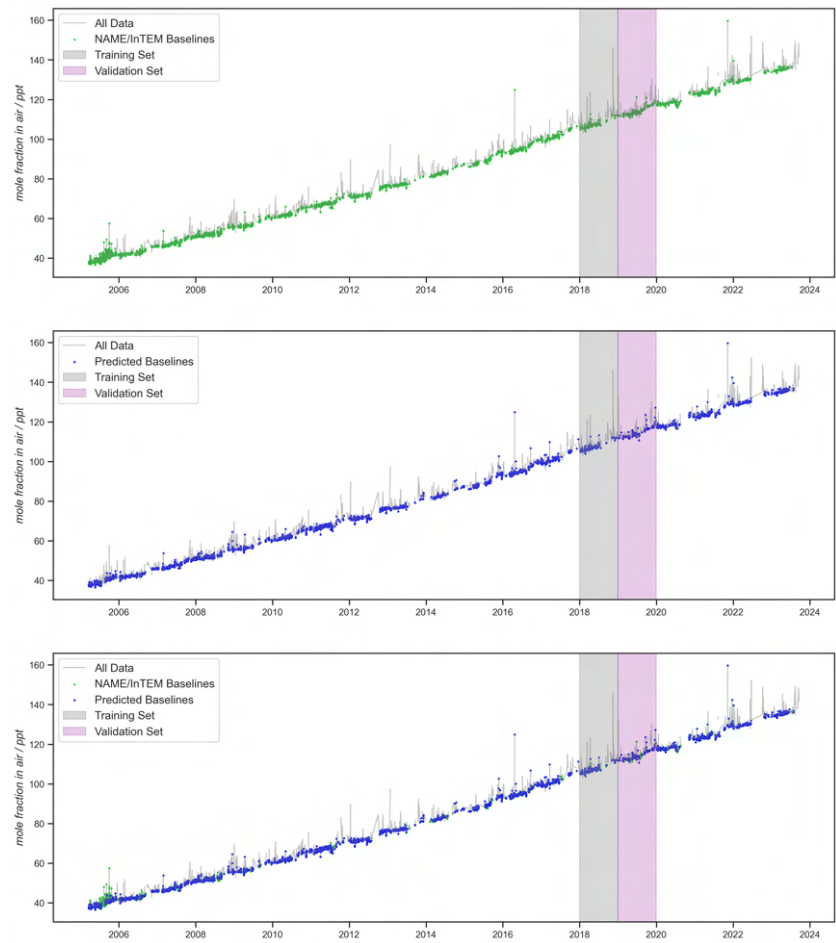


230 Monthly means

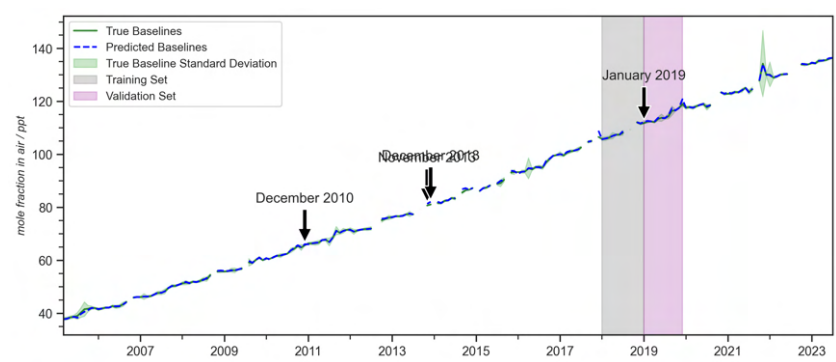


4.11.8 HFC-134a

Mole fraction time series

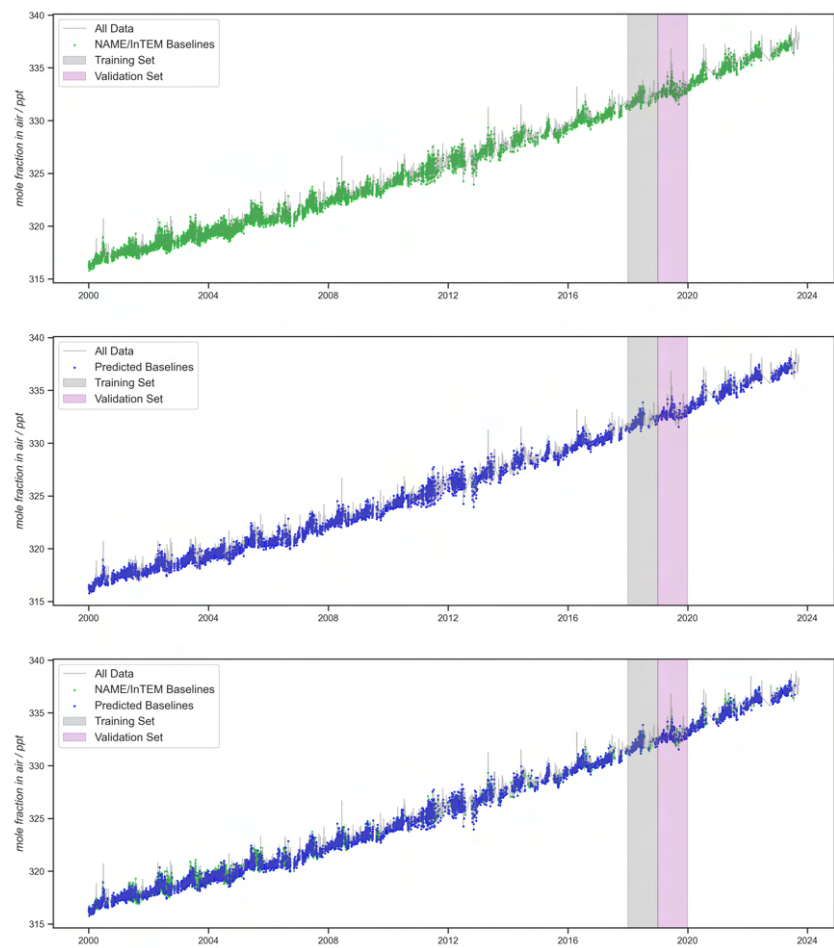


Monthly means

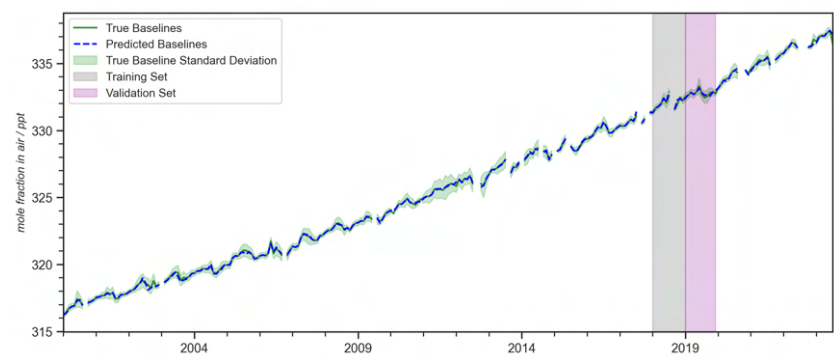


4.11.9 N₂O

235 Mole fraction time series

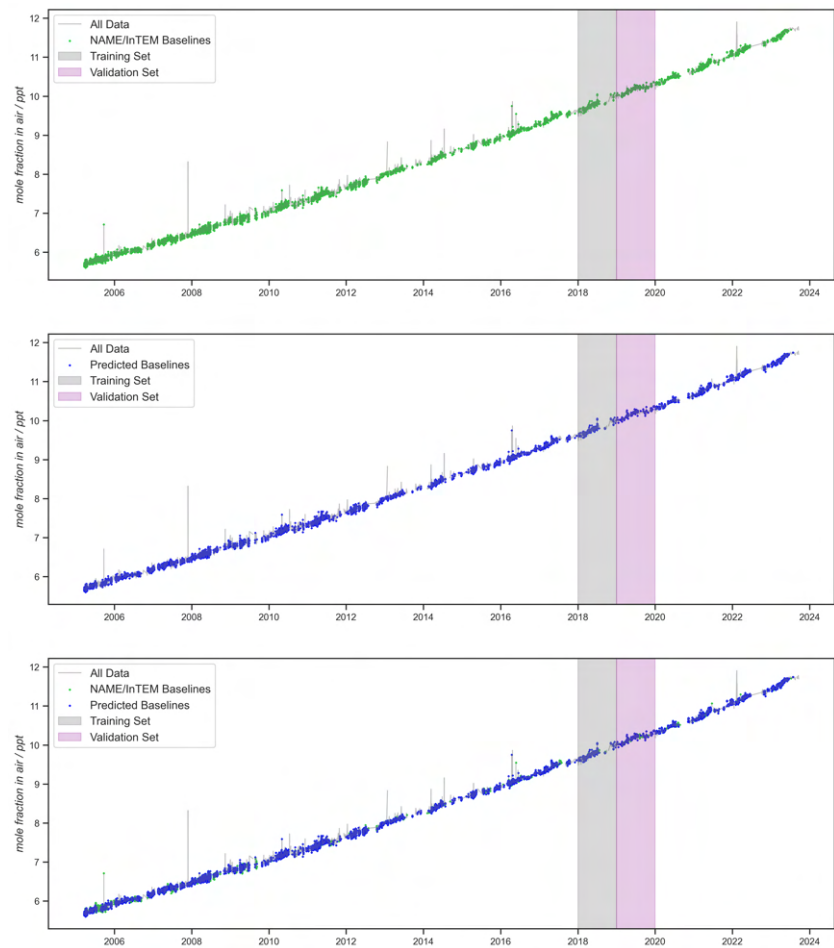


Monthly means

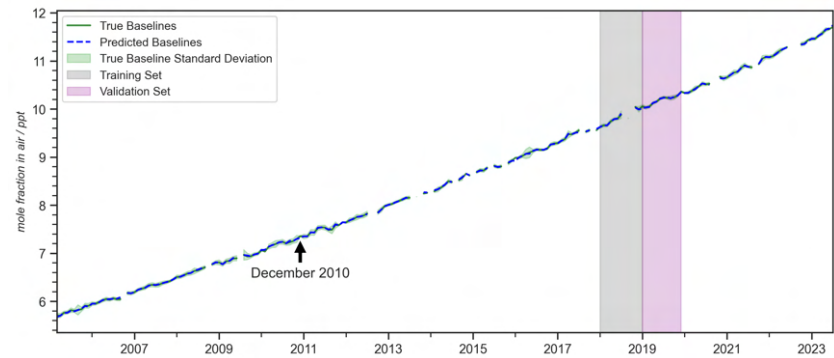


4.11.10 SF₆

Mole fraction time series

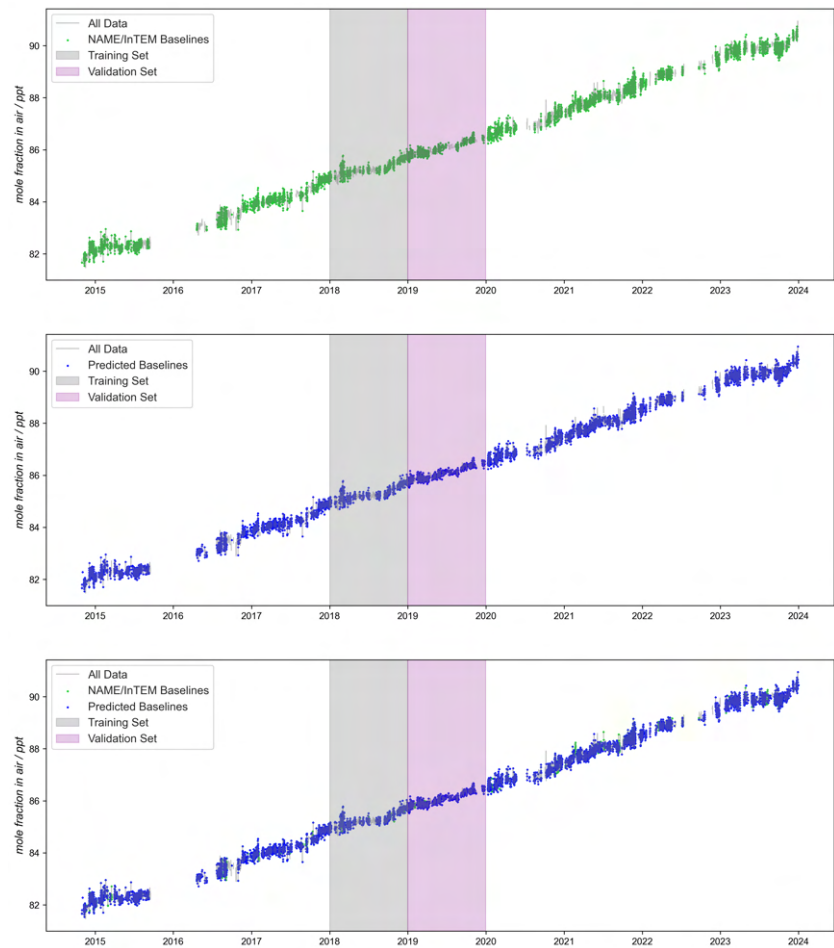


Monthly means

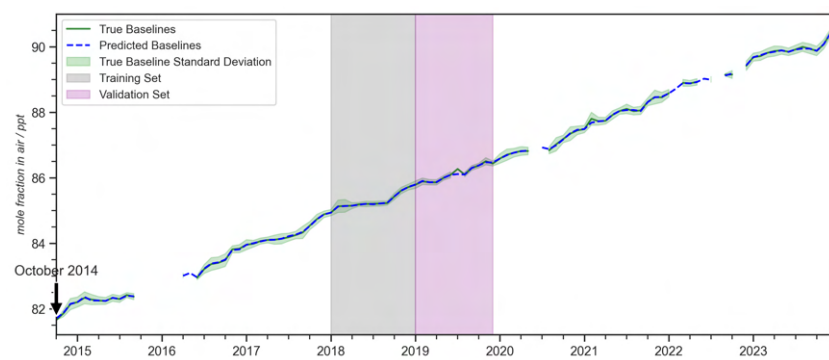


4.12.1 CF₄

Mole fraction time series

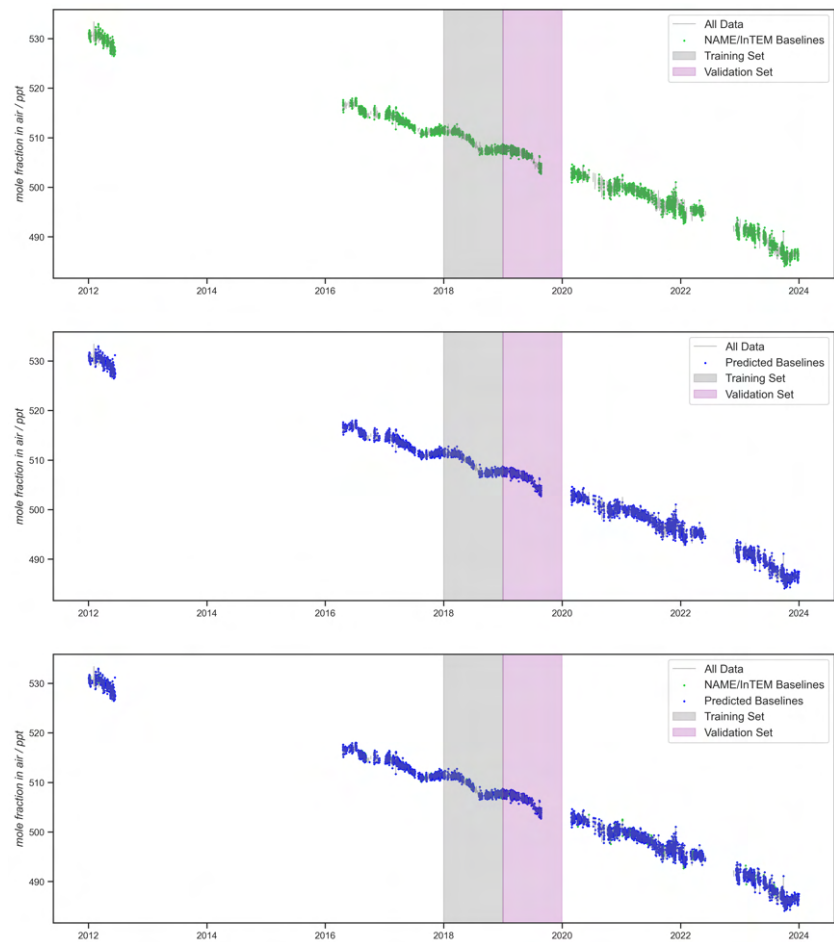


Monthly means

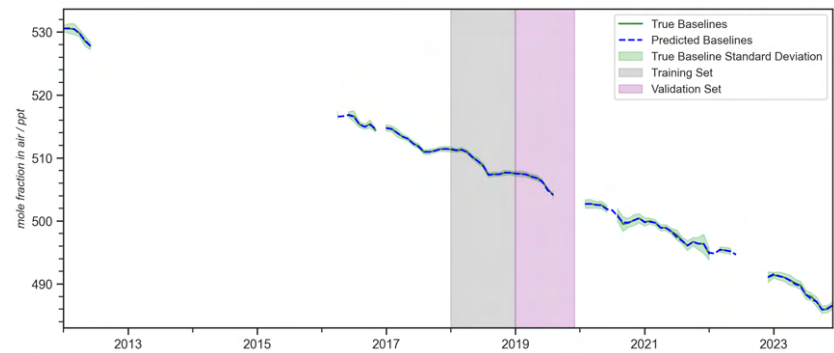


4.12.2 CFC-12

245 Mole fraction time series

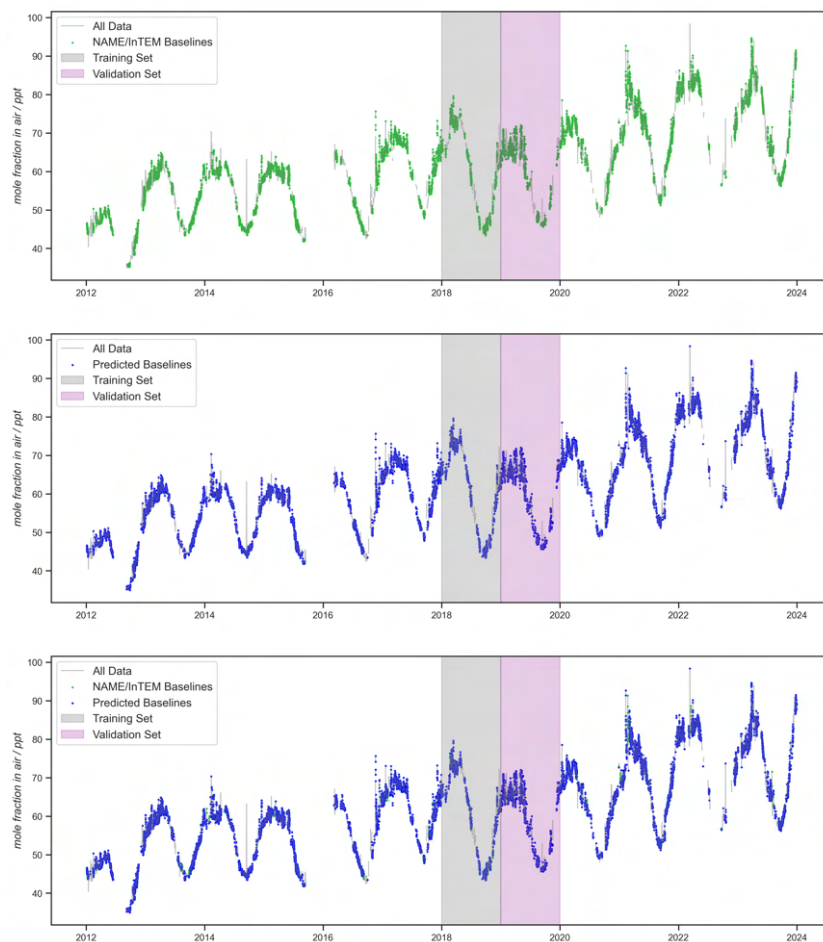


Monthly means

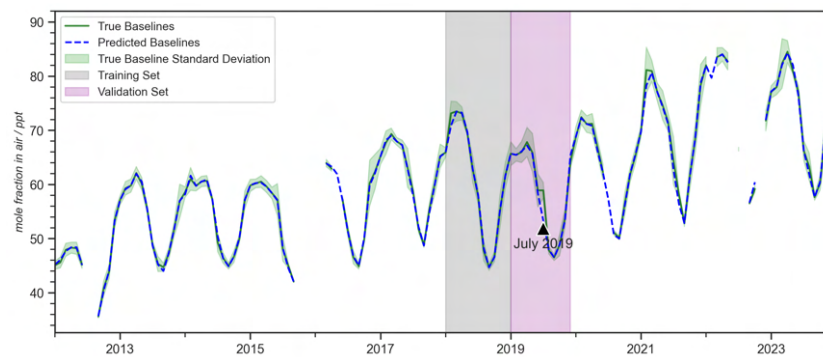


4.12.3 CH₂Cl₂

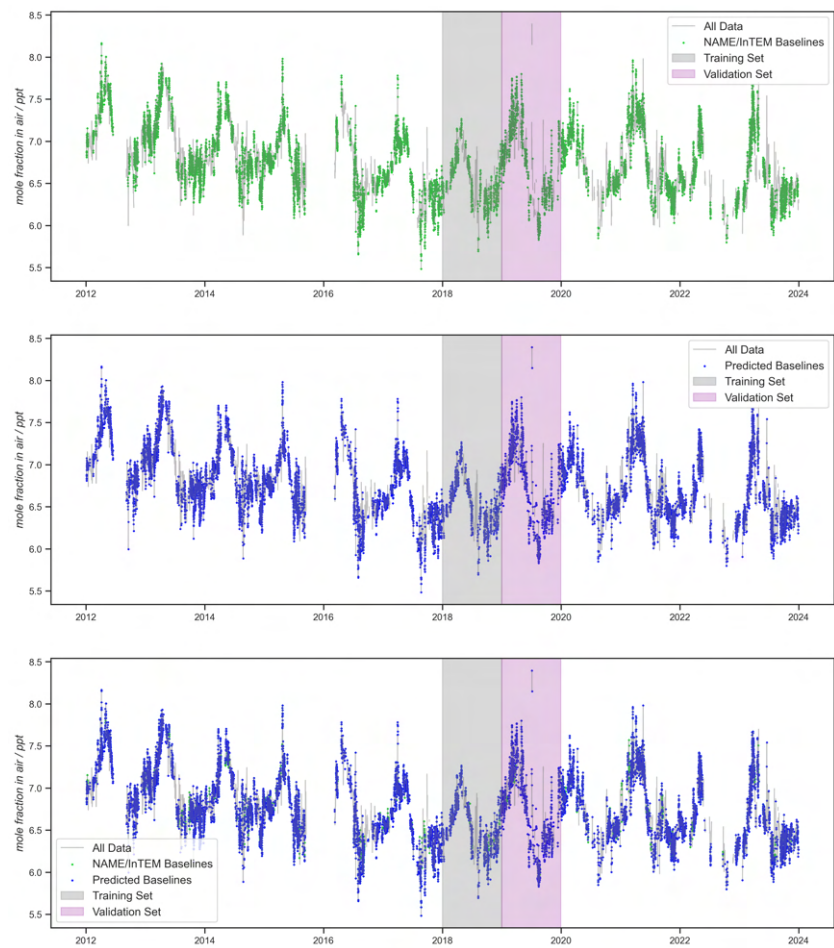
Mole fraction time series



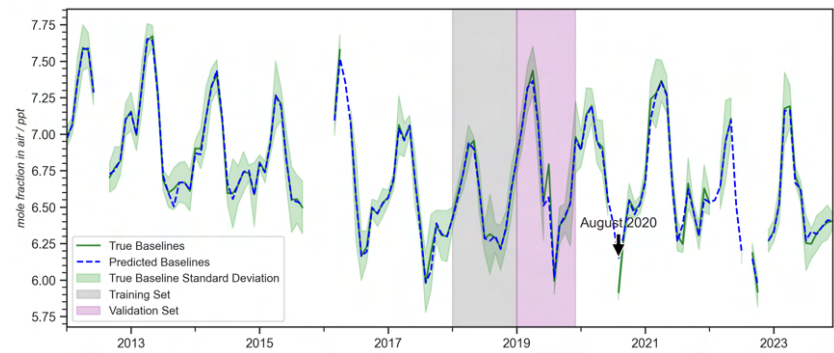
Monthly means



Mole fraction time series

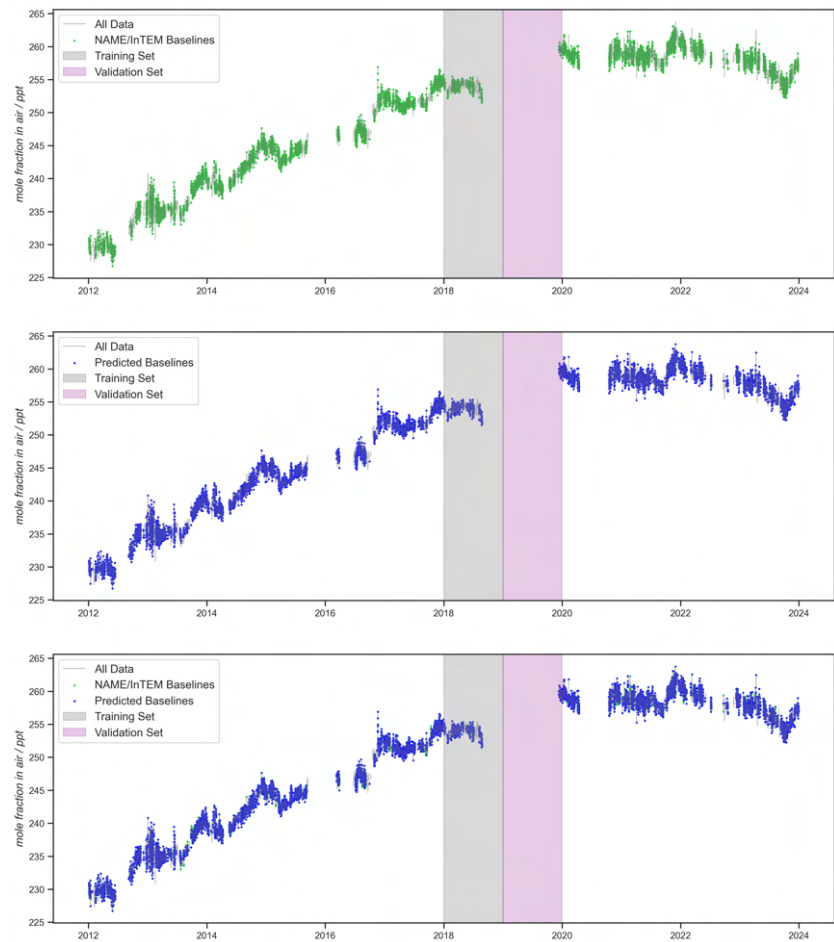


Monthly means

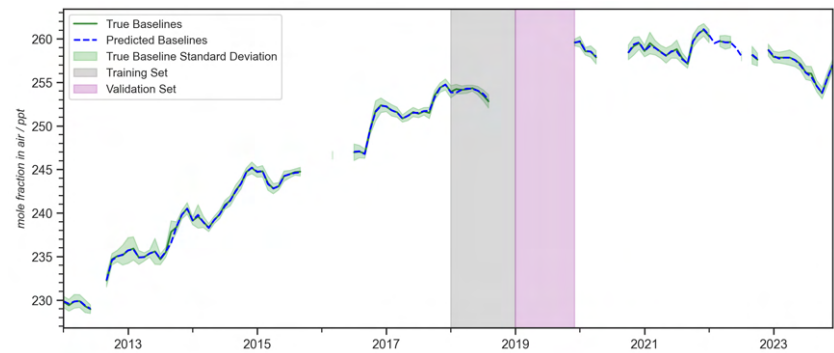


4.12.5 HCFC-22

Mole fraction time series

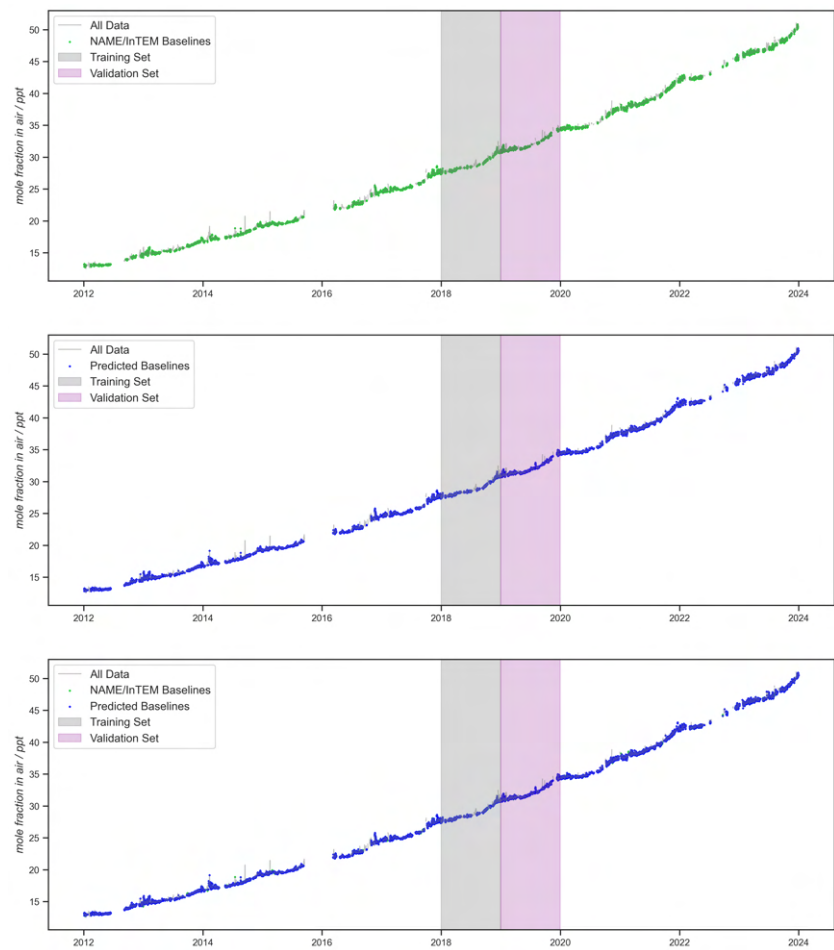


255 Monthly means

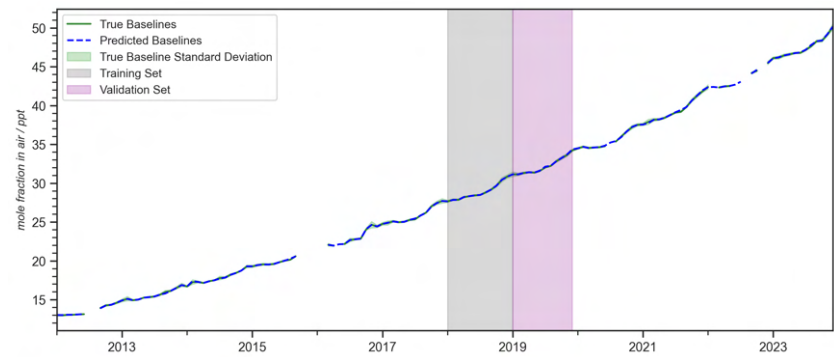


4.12.6 HFC-125

Mole fraction time series

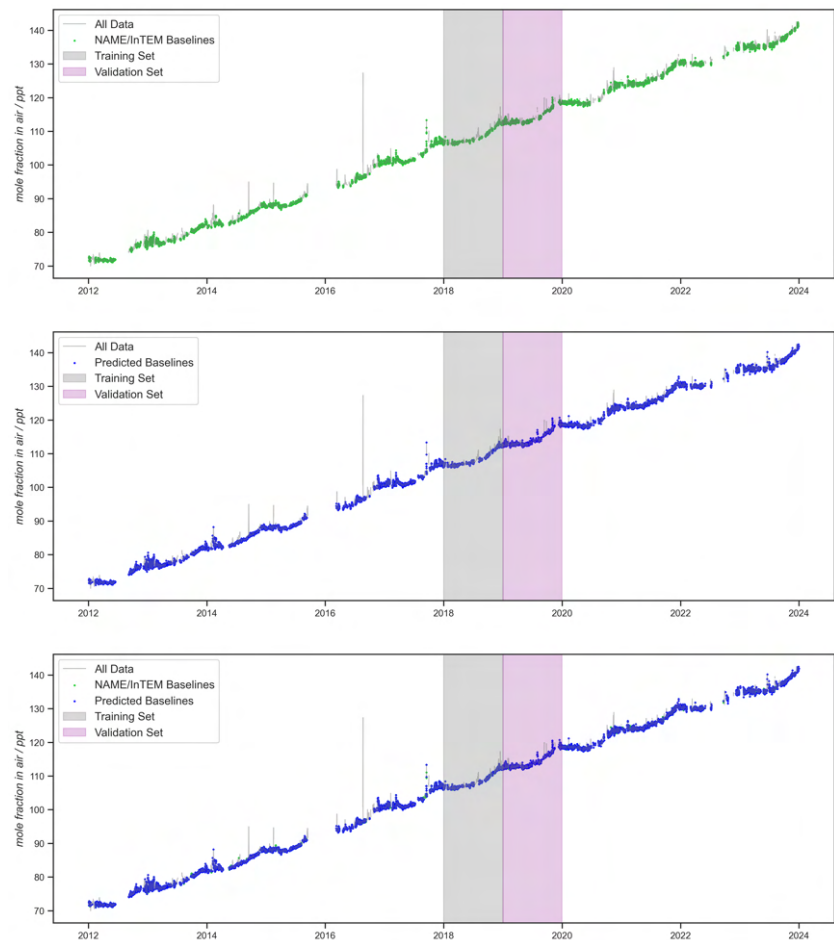


Monthly means

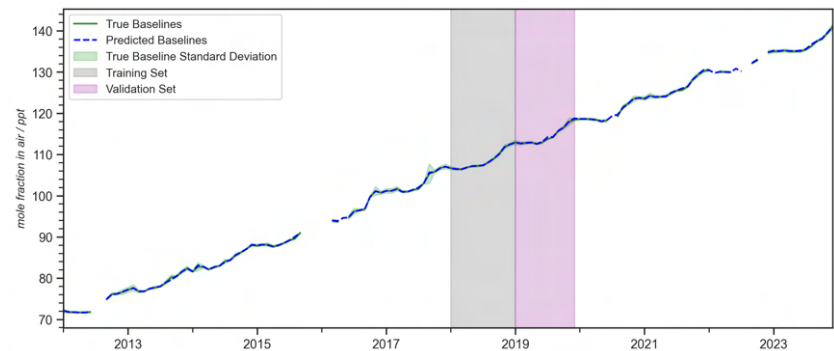


4.12.7 HFC-134a

260 Mole fraction time series

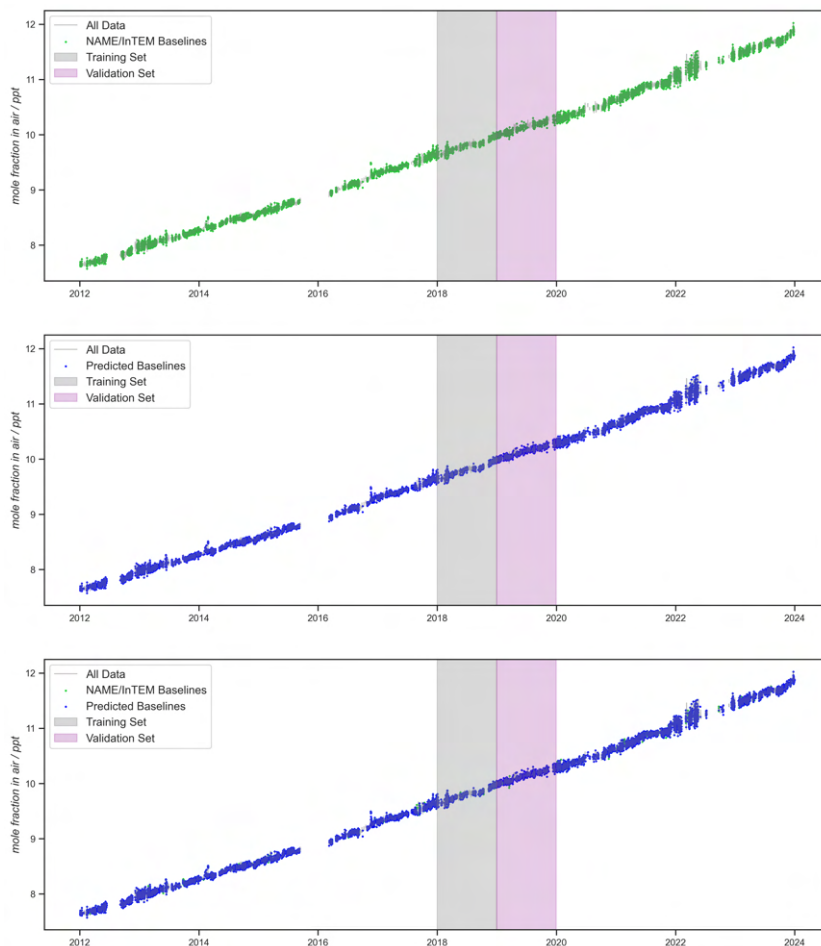


Monthly means



4.12.8 SF₆

Mole fraction time series



Monthly means

