# Geolocating Particulate Filters from the IMS Based on Machine Learning as a Means of Identifying Anomalies

B. Milbrath, A. Hagen, and C. Svinth

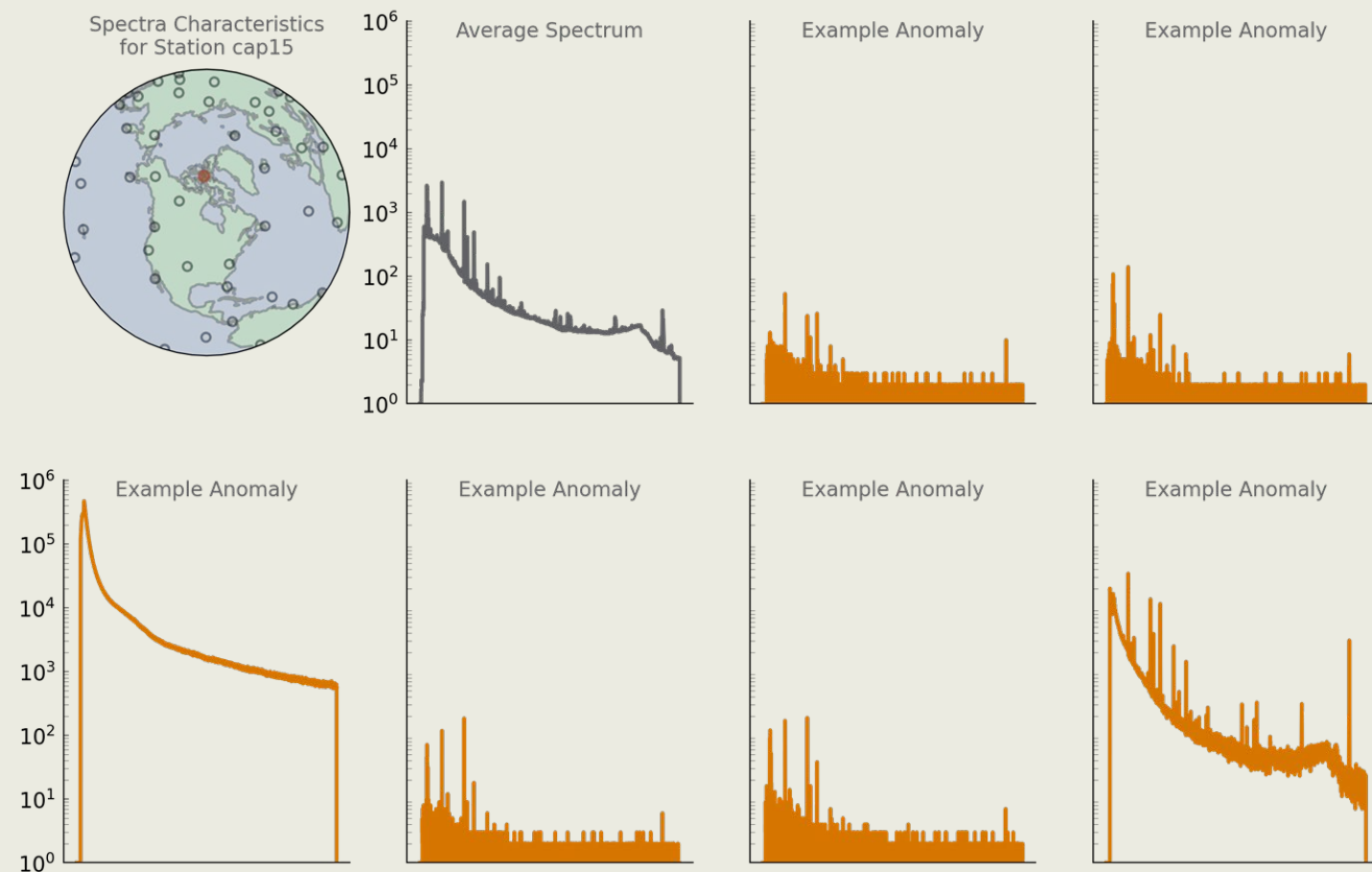Pacific Northwest National Laboratory
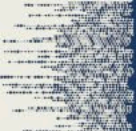
PNNL

Presentation Date: dd September 2025

B. Milbrath, A. Hagen, and C. Svinth

O3.6-188

## Motivation

- Data quality is fundamental to International Monitoring System (IMS) operation, but data corruption, mistakes, or detector degradation impact system operational effectiveness

- Automated data quality monitoring is needed

- The following case study was performed on particulate spectra data



Several example anomalies in reported spectra from station cap15

B. Milbrath, A. Hagen, and C. Svinth

# Supervised Learning

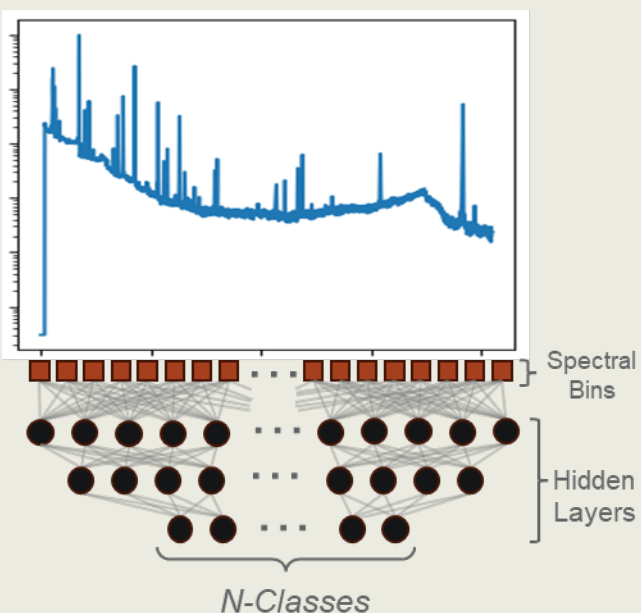We utilized supervised learning, i.e., machine learning where the network is trained on labelled data

-To do this we used a randomly-selected 80% of the data from 5 years for 29 stations (> 200k spectra)

-Verification done with the remaining 20% of the data

-Needed to make all spectra have 8192 channels

-The goal of the training is to minimize the difference between actual and predicted values
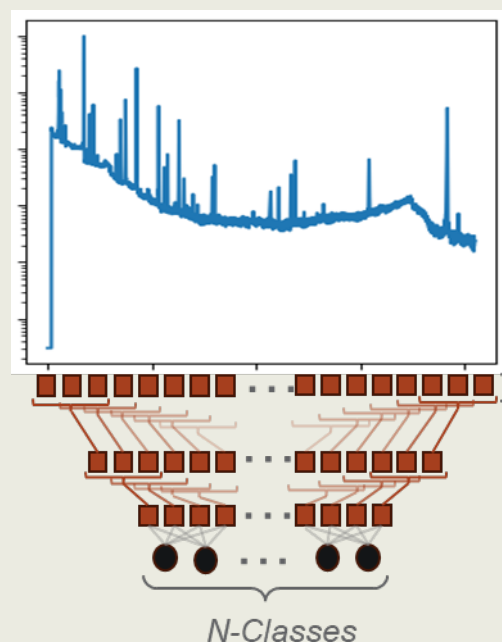
We trained neural networks to:

- Predict the location from a spectrum: network produces a location on the globe and minimizes the haversine distance from the prediction to the real location

- Predict the station from a spectrum: station classification was done by cross entropy, i.e., measuring the difference between two probability distributions
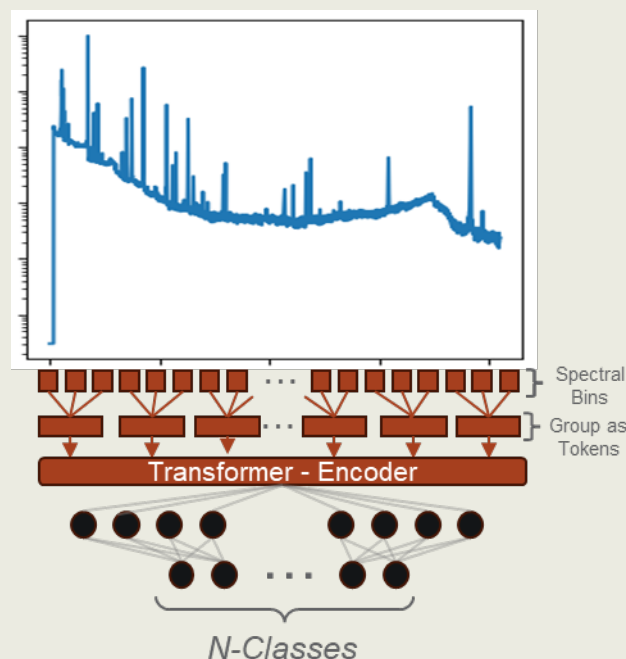
PNNL

# 3 Neural Network Types Tried



Multilayer Perceptron – fully connected neurons, i.e., information from every channel connected to others immediately

Convolutional Neural Network – less computational; good at learning features through filter optimization

B. Milbrath, A. Hagen, and C. Svinth

# 3 Neural Network Types Tried

Transformer



Transformer – newer machine learning approach; can be trained on data with different numbers of channels since the data is grouped
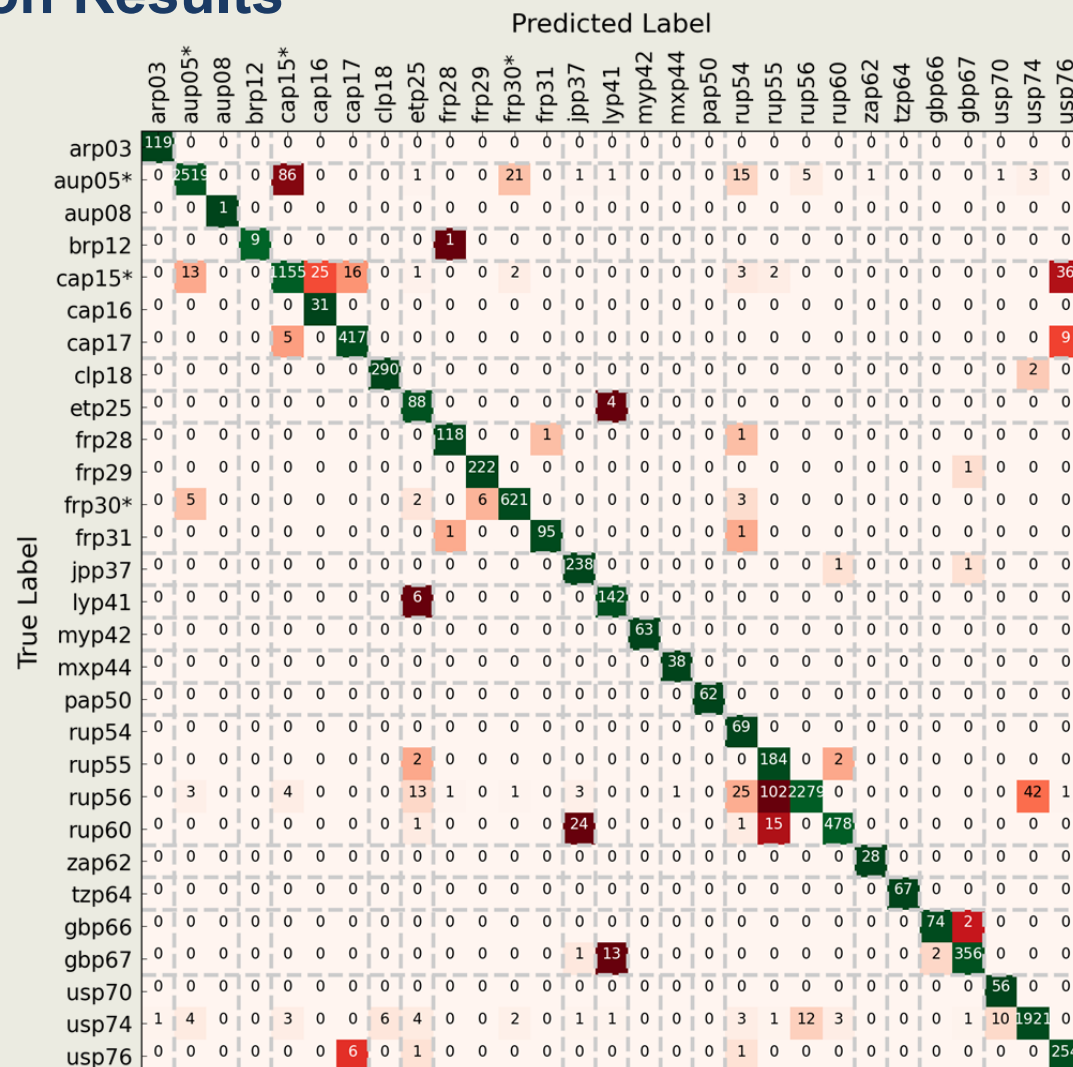
Results were about the same, no matter which neural network type used

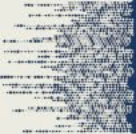>95% of location predictions within 150 km of true location (station)

97% station classification performance

# Localization Results

- Confusion matrix showing predicted station on x-axis and true station on y-axis

- White to red to purple on the off-diagonals shows the misclassifications (largest are 5-10%)

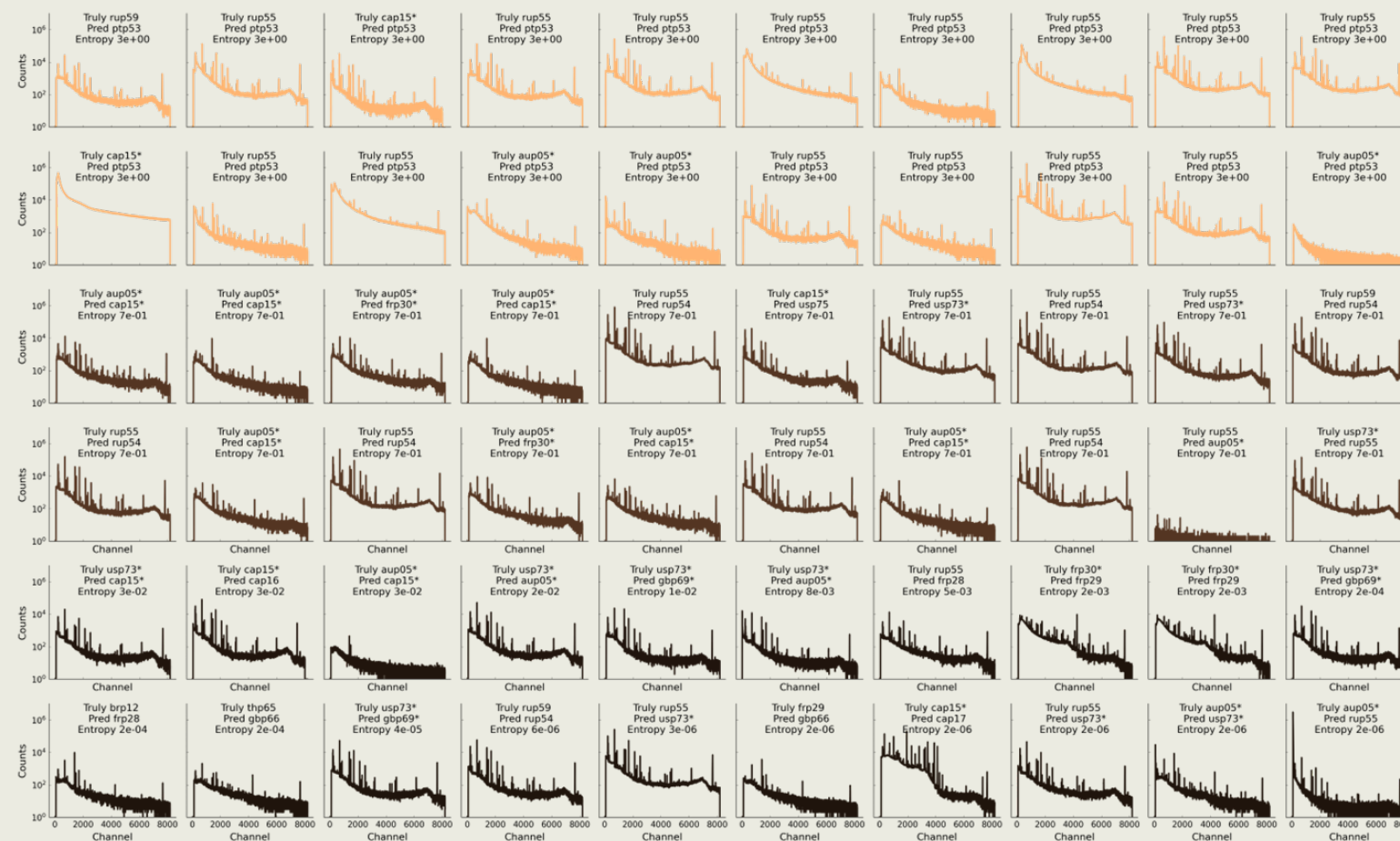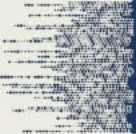-  Misidentifications are generally within a country or close region

**SnT 2025** — CTBT: SCIENCE AND TECHNOLOGY CONFERENCE
8 SEPTEMBER ONLINE DAY / 9 TO 12 SEPTEMBER AT HOFBURG PALACE, VIENNA & ONLINE

**Geolocating Particulate Filters from the IMS Based on Machine Learning as a Means of Identifying Anomalies**

O3.6-188

B. Milbrath, A. Hagen, and C. Svinth

# Mispredicted Spectra:  Electronic vs. Nuclides

- Electronic issues led to higher entropy predictions – low prediction confidence

- Nuclide differences lead to low entropy predictions – confident predictions of wrong station



Darker spectrum -> lower entropy, but all are mispredictions

PNNL

B. Milbrath, A. Hagen, and C. Svinth
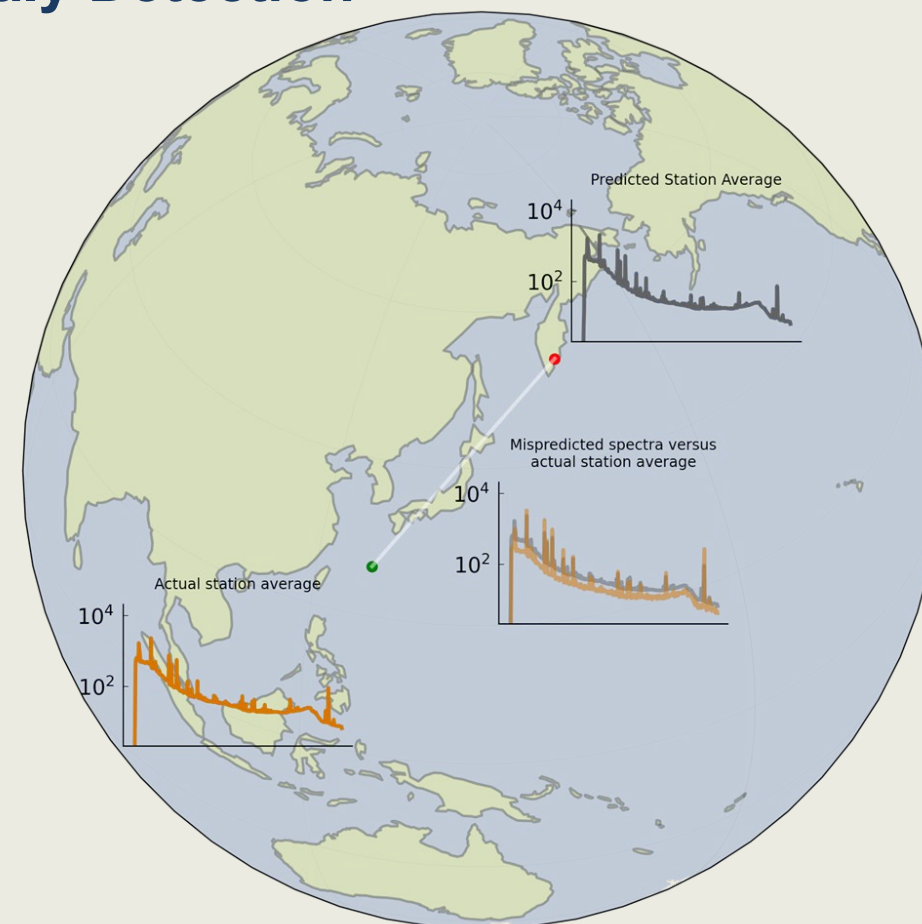
O3.6-188

# Closer Look at Mispredictions

Comparison of average mispredicted spectrum (grey) vs. predicted spectrum (orange) for stations with more than 10 mispredictions as another station
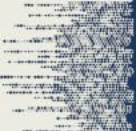
B. Milbrath, A. Hagen, and C. Svinth

# Operational Anomaly Detection

- Automated Triage
  - Data quality
  - Data corruption
  - Detector issues and changes

- Station expectations
  - Nuclides
  - Overall activity

- Connections to ATM and climatic regions could provide further context

- Classification success implies the ability to compare sets of stations or readings for trend analysis or event analysis



Automated identification of spectra more indicative of another station is followed by human triage
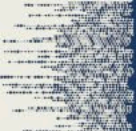
B. Milbrath, A. Hagen, and C. Svinth

## Machine Learning Operations

- Supervised learning (SL) is ideal for
  - Automation of "easy tasks"
  - Metrics
  - Effect detection
- Operationalization of SL requires
  - Alignment with mission metrics
  - Dimension reduction
  - Input monitoring and alerting
  - Automated retraining
- Risks
  - Input drift, e.g. change in expected nuclides change, detectors change
  - *Often, you must start from scratch*

### Alternatives

- Unsupervised Learning (UL)
  - Dimensions reduced and combined into a representation with appropriate invariances
  - Two sample tests or other statistical techniques performed on dimension reduced space
- Machine learning guided signature discovery
  - Techniques like the previously presented identify features which can be used without ML in the future

PNNL

**Geolocating Particulate Filters from the IMS Based on Machine Learning as a Means of Identifying Anomalies**

B. Milbrath, A. Hagen, and C. Svinth

O3.6-188

# Conclusions

- Particulate spectra reported to the IDC may have anomalous or corrupted data from
  - Electronic and maintenance state of the detector
  - Reporting issues
  - Other

- Assuring the authenticity and quality of data is critical to the IMS performance

- Particulate spectra contain information about location and time of measurement

- Supervised neural networks can predict originating station by evaluating raw spectra

- Prediction performance is 97% by station classification, 95% by location regression

- Mispredictions often indicative of anomalous spectra

- High performance implies a high information content in spectra and possibilities for later, more advanced analyses

PNNL