



Unsupervised Classification of Absorbing Aerosols with the SP2 via a Variational Autoencoder (VAE)

Aaryan Doshi¹ and Kara D. Lamb²

¹Stanford University, Palo Alto, CA, 94305

²Department of Earth and Environmental Engineering, Columbia University, New York, NY 10027

Correspondence: Aaryan Doshi (adoshi25@stanford.edu) and Kara D. Lamb (kl3231@columbia.edu)

Abstract. The Single Particle Soot Photometer (SP2) detects refractory aerosol particle mass on a single-particle basis via laser-induced incandescence (L-II). While the SP2 has traditionally been used to quantify black carbon aerosol mass in the atmosphere, the instrument is increasingly being used to detect and quantify other types of absorbing aerosols, such as mineral dust or anthropogenically-sourced iron oxide aerosols. Quantifying the mass loadings and emission sources of absorbing aerosols in the atmosphere is important for understanding their role in the climate cycle. Supervised machine learning algorithms have shown potential to classify different types of aerosols from L-II signals, but these methods are sensitive to instrument configuration and require training datasets generated from laboratory samples, which do not generalize well to ambient atmospheric aerosols. Here we explore the effectiveness of an unsupervised deep learning method, a variational autoencoder (VAE), applied directly to L-II signals from the SP2 in order to classify different types of absorbing aerosols. The VAE compresses L-II signals into a bottleneck latent representation and reconstructs an output as similar as possible to the input signal, thereby reducing dimensionality. We apply this approach to a dataset comprised of laboratory samples of materials that show detectable incandescence in the SP2, including fullerene soot (as a proxy for black carbon), coated fullerene soot, coal fly ash, mineral dust, volcanic ash, hematite, and magnetite. We explore optimal latent representations of L-II signals to maximize separability of different aerosol classes by varying the size of the latent representation, and find that a latent representation of 3 allows us to capture the majority of the information in the L-II signals relevant for identifying different types of absorbing aerosols. We demonstrate that unsupervised machine learning is a promising method for identifying distinct populations of aerosols detected by the SP2.

1 Introduction

Atmospheric aerosols originate from a multitude of sources: primarily from natural sources such as wind-blown dust from deserts, sea salt from the oceans, and smoke emitted from forest fires, and from anthropogenic sources such as the combustion of fossil fuels. Detection of these aerosols is essential to understanding the impact they have on the climate. While some aerosols such as sulfate aerosols from volcanoes can have a cooling effect by blocking out sunlight, other types of aerosols absorb light from the sun, heating up the atmosphere locally and substantially contributing to the overall warming of climate (Bond et al., 2013; Baumgardner et al., 2012). Absorbing aerosols such as black carbon, brown carbon, and mineral dust are



25 important short-lived climate forcers, with significant direct climate radiative effects. Quantifying the atmospheric abundance of these aerosols requires *in situ* observations in order to determine their emissions, sources, and lifetime in the atmosphere.

The Single Particle Soot Photometer (SP2), shown in Figure 1, is the state-of-the-art instrument for detecting refractory black carbon (rBC) in the atmosphere (Stephens et al., 2003). For the past two decades, the SP2 has been used in numerous ground-based and airborne field studies to measure the atmospheric abundance of rBC (e.g. Schwarz et al., 2006; Moteki and Kondo, 2010; Lamb et al., 2018), an aerosol sourced from incomplete combustion that has important implications for climate (Bond et al., 2013). Increasingly, the SP2 is being used to detect a more diverse spectrum of light-absorbing refractory aerosols like mineral dust and anthropogenically sourced iron-oxide aerosols (Moteki et al., 2017; Liu et al., 2018; Lamb, 2019).

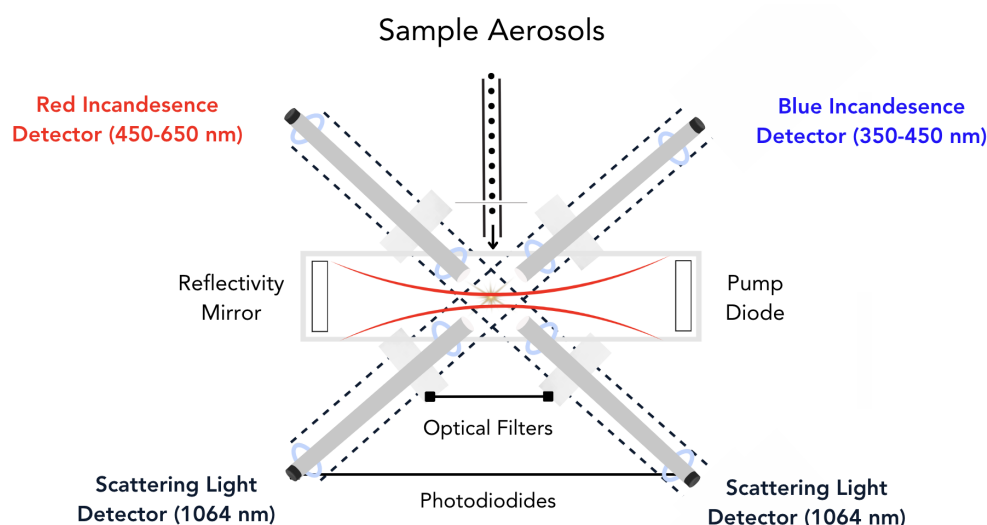


Figure 1. A schematic of the SP2 instrument. Figure adapted from Schwarz et al. (2006).

The SP2 uses laser-induced incandescence (L-II) to quantify refractory aerosol mass on a single particle basis (Stephens et al., 2003). The SP2 samples aerosols in the sub-micron range by pulling particles into the cavity of an ND-YAG laser (1064 nm). As particles pass through the laser beam, if they have a sufficient absorption cross-section at 1064 nm, they will heat up and incandescence. Here we use observations from the NOAA SP2 instrument (Schwarz et al., 2006, 2010). In its typical configuration, the NOAA SP2 acquires signals on 4 detectors with a 5 MHz acquisition rate as the aerosols pass through the center of the laser beam. Two channels detect light scattered by particles as they pass through the center of the ND-YAG laser using an avalanche photo-diode, and two channels detect light emitted by particles during incandescence using photomultiplier tubes (PMT) that measure visible light in two spectral bands: a narrow band PMT with peak sensitivity at 420 nm (350-450 nm) that we refer to as the “blue” incandescent channel, and a broadband PMT with peak sensitivity at 630 nm (450-650 nm) that we refer to as the “red” incandescent channel. One of the scattering channels has a position-sensitive detector and is used to determine the position of the particle relative to the center of the laser beam, as has been described in detail in Gao et al. (2007). The position sensitive detector is used to determine the center point of the laser in order to derive information about the coating



45 state of rBC particles by using the leading-edge-only fitting method and assuming Mie core-shell theory (Gao et al., 2007). The L-II signal associated with each aerosol detected by the SP2 therefore consists of time series from these 4 detection channels, which provide information about how the particle scatters and emits light as it passes through and (potentially) evaporates in the laser beam.

Supervised machine learning, involving the training of algorithms on a labeled dataset, has been used in the past to classify
50 different types of absorbing aerosols with detectable incandescence in the SP2 (Lamb, 2019). While supervised machine learning methods can classify aerosols based on features derived from L-II signals, the algorithms need to first be trained on labeled data sets. Lamb (2019) previously used observations of laboratory proxies for typical atmospheric aerosols detectable by the SP2 to create a labeled data set. While supervised machine learning performs well in classifying laboratory data sets by type, it has limited ability to generalize to ambient atmospheric aerosols, especially in cases where laboratory proxies are
55 not readily available or when aerosol populations that have not been previously identified in past data sets are measured during airborne field campaigns.

On the other hand, unsupervised machine learning algorithms discern inherent patterns and correlations in data sets, without requiring predefined categories. Unsupervised machine learning has not previously been applied to the problem of classifying L-II signals. Here we explore how unsupervised machine learning can be applied to L-II signals from the SP2, with the goal
60 of identifying different populations of aerosols based on the information in their L-II signals alone. We explore how these unsupervised methods can provide insights into the variability of aerosols of different types based on their L-II response in the SP2. Our analysis also provides insight into the amount of independent information that can be gained from L-II signals in terms of identifying the composition of refractory aerosols that reach detectable incandescence in the SP2.

To demonstrate this method, we focus on the application of unsupervised machine learning to observations of laboratory
65 proxies for several types of anthropogenic and natural aerosols that reach detectable incandescence in the SP2 (Lamb, 2019). We briefly describe these data sets, data pre-processing, and the unsupervised machine learning algorithm in Section 2. We then discuss the application of the VAE to the laboratory samples in Section 3. In Section 4 we discuss how this approach can be used to identify outliers in ambient populations, and in Section 5 we explore how this method can be used to improve classification of different aerosols that the SP2 is sensitive to. Finally, in Section 6, we discuss the potential for machine learning to improve
70 SP2 data analysis and interpretation.

2 Methods

2.1 Dataset

To investigate how effectively unsupervised machine learning can be used to differentiate different types of aerosols detected by the SP2, we use a labeled data set of L-II time series (Lamb, 2025) that was previously described in detail in Lamb (2019). This
75 dataset is comprised of L-II signals obtained from measuring laboratory proxies for aerosols typically found in the atmosphere that reach detectable incandescence in the SP2. The data set includes examples of observations for 7 classes of aerosols: Fullerene Soot (FS), Fullerene Soot coated with glycerol (FS+glyc), Clifty Fly Ash (CFA), Arizona Test Dust (ATD), Volcanic



Table 1. Overview of laboratory data sets.

Label	Class	Total	% with Detectable Incandescence
0	FS	20004	98.28
1	FS+glyc	20018	81.04
2	CFA	20009	8.53
3	ATD	20001	23.93
4	VA	20005	27.00
5	Fe ₂ O ₃	20008	91.85
6	Fe ₃ O ₄	20037	98.21

Ash (VA), Iron (III) Oxide (Fe₂O₃), and Iron (IV) Oxide (Fe₃O₄). The total number of aerosols of each class measured is given in Table 1.

80 This laboratory data set was developed to create a relatively balanced data set that includes examples of aerosols that the SP2 might observe in the atmosphere. In ambient conditions, the majority of aerosols that the SP2 observes are rBC particles. Fullerene soot is a laboratory proxy with a similar response in the SP2 as rBC. BC observed in the atmosphere is typically coated with non-absorbing materials, and we use the FS+glyc samples as examples of typical L-II signals of thinly coated rBC particles. Several studies have demonstrated that SP2s that have been modified to provide greater spectral contrast between
 85 their narrow and broadband detectors (such as the NOAA SP2 used in this study) can also detect iron oxide aerosols associated with anthropogenic combustion sources with high efficiency: magnetite (Fe₃O₄) can be detected with nearly 100% efficiency under typical conditions, while the detection of hematite (Fe₂O₃) is lower and size-dependent (Yoshida et al., 2016). The two iron oxide powders (Fe₂O₃ and Fe₃O₄), have a similar response in the SP2 as anthropogenically sourced iron oxide aerosols from combustion sources (Moteki et al., 2017; Lamb, 2019; Lamb et al., 2021), which we refer to as FeO_x following past
 90 literature. In addition, the SP2 also detects measurable incandescence in a small fraction of aerosols with metallic inclusions, such as mineral dust, coal fly ash, and volcanic ash (Heimerl et al., 2012; Lamb, 2019). Here we use Arizona Test Dust as a laboratory proxy for mineral dust, and the coal fly ash is Clifty-F. The volcanic ash was collected on the ground in Iceland from the Eyjafjallajökull Volcano.

2.2 Pre-processing L-II Time Series from the SP2

95 In this study, we focus on applying unsupervised machine learning to L-II signals from the SP2 instrument. In a previous study using supervised machine learning to classify aerosols detected by the SP2, significant feature engineering was used to derive specific, interpretable features from the L-II time series (Lamb, 2019). By contrast, as input to our unsupervised machine learning algorithm, we use the unprocessed L-II signals, which are 80 μ s time series consisting of 400 points ($dt = 0.2 \mu$ s, see Figure 2 for examples), under the assumption that the non-linear deep learning method will learn higher order features directly
 100 from the unprocessed L-II time series, that will provide insights into the variability and distinguishability of the observed aerosol particles.



For each of the 4 detection channels, we define a feature matrix \mathbf{X}_i , where $\mathbf{X}_i \in \mathbb{R}^{N \times t}$ is an $N \times t$ matrix. N is the total number of L-II signals (the total number of observed aerosol particles), $t = 400$ corresponds to the number of time points in each signal, and $i \in [0, 3]$ corresponds to the detection channel in the SP2 instrument. The 0th channel corresponds to the scattering channel, the 1st channel corresponds to the “blue” incandescent channel, the 2nd channel corresponds to the “red” incandescent channel, and the 3rd channel corresponds to the position sensitive detector.

The typical approach to derive information from L-II signals from the SP2 is to find the maximum values for the scattering and incandescent channels, as these are proportional to the optical size and refractory mass of an aerosol, respectively. We define the maximum of channel 0 as $S_{max} = \max(\mathbf{X}_0)$, and the maximum of channel 1 as $I_{max} = \max(\mathbf{X}_1)$. In addition, the “color temperature ratio”, is defined as

$$CR = \frac{\max(\mathbf{X}_1)}{\max(\mathbf{X}_2)}, \quad (1)$$

the ratio between the peaks of the blue and red incandescent signals. CR is proportional to the blackbody temperature of the aerosol as it incandesces in the laser beam. Here the gains on the blue and red detectors have been chosen such that the $CR \approx 1$ for rBC (corresponding to a characteristic blackbody temperature of 4320 K) and $CR \approx 0.7$ for FeO_x (corresponding to 3300 K). In practice, there is a significant amount of variability in CR across the population of aerosols of each class detected by the SP2. For further details of typical SP2 analysis, see discussion in Schwarz et al. (2006, 2010); Lamb (2019).

To first order, FeO_x signals can be differentiated from rBC signals in the SP2 by differences between their blackbody temperature (CR) and their incandescent mass (proportional to the peak of the incandescent channel, I_{max}). However, CR and I_{max} do not provide complete separation between the two classes, particularly when I_{max} is small (i.e. for less massive particles) (Lamb, 2019). FeO_x also demonstrates a less skewed incandescent peak in the SP2 than rBC signals, likely due to the metallic aerosols melting in the SP2 laser beam as they are heated to incandescence (Adachi et al., 2016; Lamb, 2019). In addition, other types of aerosols detected by the SP2 such as ATD, VA, and CFA exhibit a broad range of I_{max} and CR values, due to the presence of metallic inclusions with a variety of chemical compositions. The maximum value of the scattering channel, S_{max} , is proportional to the total optical size of the aerosol particle (except in cases when particles are large and the scattering channel is saturated). ATD, VA, and CFA generally also demonstrate significant scattering in the laser beam after the main incandescent peak, due to incomplete evaporation of these aerosols in the SP2. Because these particles are generally larger than typical rBC particles, the scattering channel is more likely to be saturated for these aerosol types. Coated rBC and FeO_x particles can be identified from the Ch. 0 time series due to an initial peak in the scattering signal as the coating evaporates from the particle, followed by a second peak when the refractory portion of the particle evaporates in the laser beam.

In this study we focus on particles that have detectable incandescence in the SP2. Therefore, we first remove any L-II signals where I_{max} is close to the Ch. 0 signal baseline ($I_{max} < 3$). The majority of the rBC and FeO_x aerosols have detectable incandescence. However, only a fraction of the ATD, CFA, and VA aerosols demonstrate detectable incandescence in the SP2 (Table 1), likely because only a fraction of the particles in these aerosol populations have sufficient metallic inclusions (Lamb, 2019).



135 Pre-processing data is a typical first step for applying deep-learning algorithms to data sets, as methods work best when the input features are normalized between 0 and 1 and normally distributed. In our analysis, we tested several potential methods for pre-processing Ch. 0 and Ch. 1 time series, which we delineate here:

1. *Division by the maximum of each channel across all samples in the training dataset.*
2. *Normalization of each channel by the minimum and maximum across all samples in the data set*
- 140 3. *Normalization of the channel by the minimum and maximum of each individual sample in the data set*
4. *Relative scaling for each sample using logarithmic normalization.*

We found that the third approach demonstrated the most promise in terms of separability of classes within the latent space learned by the VAE analysis that we describe in the rest of this paper. To pre-process the raw L-II signals for unsupervised machine learning, we therefore normalized Ch. 0 and Ch. 1 such that the time series for the scattering and blue incandescent
145 channels are normalized between 0 and 1. That is, the input feature vector for our algorithm is

$$\mathbf{X}_i^{\text{scaled}} = \frac{\mathbf{X}_i - \min(\mathbf{X}_i)}{\max(\mathbf{X}_i) - \min(\mathbf{X}_i)} \quad (2)$$

where $i \in [0, 1]$. The pre-processing method that we choose impacts the meaning of latent variables that are learned by the VAE. In this case, the normalization approach that we have chosen means that we focus on learning compressed latent representations that describe the shape of the signals for Ch. 0 and Ch. 1, under the assumption that the shape of the L-II signal
150 alone (without information about magnitude) can provide information that can be used to differentiate between different types of aerosol particles that are detected by the SP2.

2.3 Dimensionality Reduction of L-II Signals with a Variational Autoencoder

After preprocessing the raw L-II signals for Ch. 0 and Ch. 1, the data sets are randomly split into training, validation, and test data sets based on the aerosol sample number. We use 50% (140,082 samples) for training, and 25% (70041) for validation and
155 25% (70042) for testing our unsupervised machine learning approach.

To apply unsupervised machine learning to the L-II signals, we use a variational auto-encoder (VAE). A VAE is a type of generative machine learning model that is designed to generate new data that is similar to the data that it is trained on (Kingma and Welling, 2022). It does this by learning a non-linear mapping from a higher-dimensional feature space to a lower dimensional latent space representation, which can then be sampled to generate new data points. VAEs are widely used in
160 machine learning for tasks like image generation and anomaly detection (Wei et al., 2020).

A VAE consists of two neural networks with trainable weights, which are the encoder and the decoder models (Figure 2, bottom panel). The training of a VAE involves optimizing the parameters of both the encoder and decoder in order to maximize the reconstruction of the original higher-dimensional input from its lower-dimensional latent space representation, while also ensuring that the learned latent space representation is smooth and continuous. The encoder takes an input signal

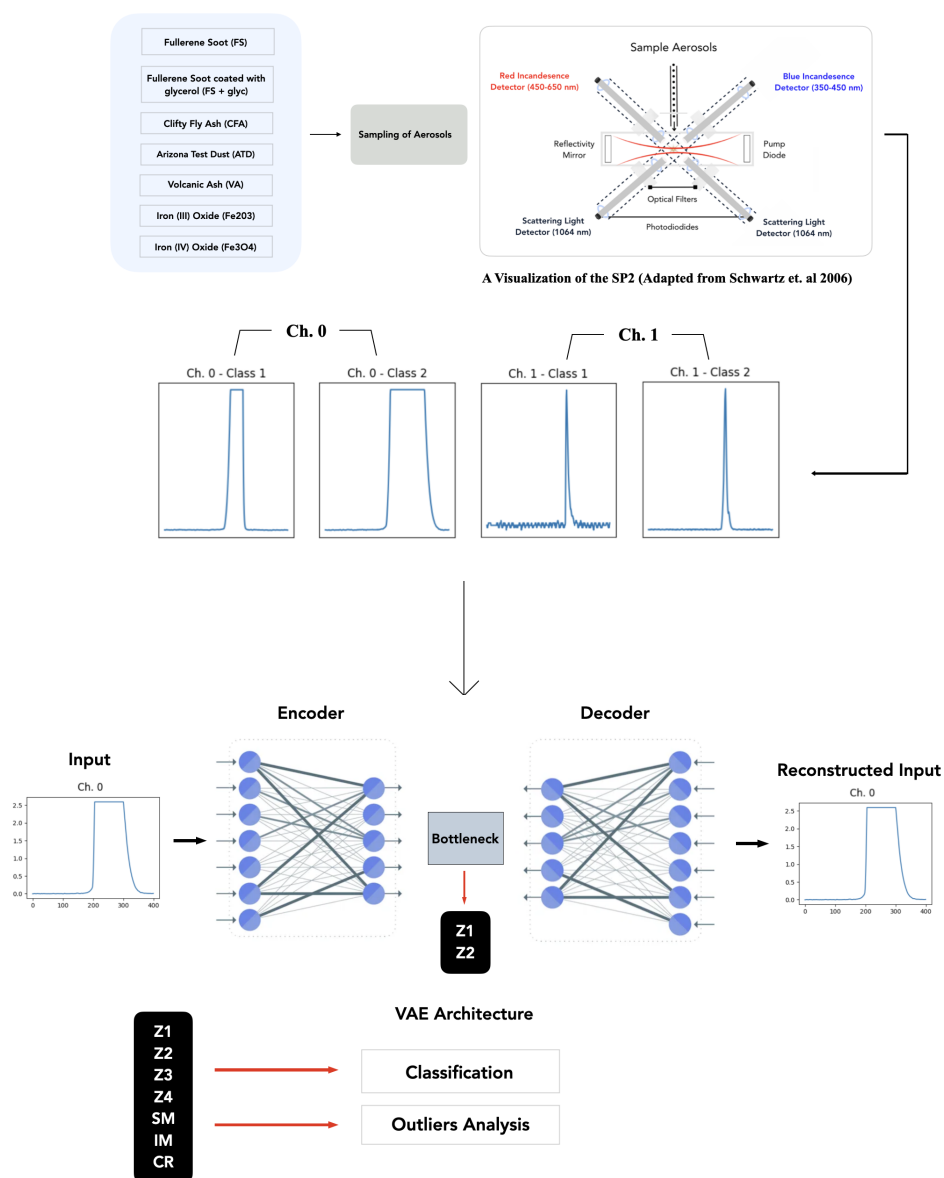


Figure 2. An overview of the unsupervised machine learning approach applied to observations from the SP2. First, aerosols are sampled by the SP2 and then after pre-processing, a variational autoencoder is trained to learn a lower dimensional latent representation of the L-II signals from each detection channel. We then explore several downstream tasks, using the latent space representations learned from the L-II signals.

165 and encodes it into a lower dimensional latent representation (Figure 2). The action of the encoder can be represented as a function, $q(z|x, \theta_{enc})$, where x is the input signal, z is its latent representation, and θ_{enc} are the weights of the encoder neural network. The encoder outputs parameters to a probability distribution, assumed to be Gaussian, which are the mean μ and



variance $\log(\sigma^2)$. Meanwhile, the decoder takes a point z in the latent space and reconstructs the input \hat{x} . The decoder defines a probability distribution over the possible outputs given a latent point. For this reason, the decoder can be represented as
 170 $p(\hat{x}|z, \theta_{dec})$ where θ_{dec} are the weights of the decoder neural network, and \hat{x} is the reconstructed signal (Kingma and Welling, 2022).

In order to learn the weights θ_{enc} and θ_{dec} of the encoder and decoder networks, we minimize the loss function for a VAE, which is the sum of the reconstruction loss and the similarity loss (Kingma and Welling, 2022). The reconstruction loss ensures that the decoded samples match the original inputs, while the similarity loss ensures that the learned latent representation is
 175 smoothly varying.

For reconstruction loss, we use binary cross-entropy loss.

$$L_{\text{reconstruction}} = -q(z|x) [\log p(x|z)] \quad (3)$$

We also explored using mean squared error loss for the reconstruction loss and found that this did not make a significant difference in our analysis of the L-II signals.

180 The similarity loss ensures that the distribution of latent variables (z) stays close to a prior distribution, which is assumed to be a standard normal distribution (Kingma and Welling, 2022). This term acts as a regularizer and effectively ensures that the learned latent space is smoothly varying:

$$L_{\text{similarity}} = D_{\text{KL}}(q(z|x) || p(z)), \quad (4)$$

where D_{KL} is the Kullback-Leibler divergence between the encoder's distribution $q(z|x)$ and the prior distribution $p(z)$. The
 185 KL-divergence measures the distance between two data distributions, and is defined as,

$$D_{\text{KL}}(p(x) || q(x)) = \int -p(x) \ln \left(\frac{p(x)}{q(x)} \right) dx \quad (5)$$

The total loss L_{total} for the VAE can then be computed by summing the reconstruction loss and similarity loss:

$$L_{\text{total}} = L_{\text{reconstruction}} + L_{\text{KL}} \quad (6)$$

Here we use the VAE algorithm implemented in the pyroVED library (Ziatdinov; Biswas et al., 2023), which is built on top
 190 of the Pytorch deep learning library and the Pyro probabilistic programming language (Bingham et al., 2018). The pyroVED library minimizes Eq. 6 using stochastic variational inference, using the Adam optimizer (Kingma and Ba, 2017).

For the L-II signals from the SP2, we independently train two VAE's on the normalized Ch. 0 and Ch. 1 signals, respectively. Here, we are interested in extracting information about the shape of the L-II signals from the scattering and incandescent channels, under the assumption that the shape of the signals can provide information about the type and characteristics of the
 195 aerosol that was measured in the SP2. The latent variables z effectively characterize the shape of the normalized L-II signals, and we train two VAEs in order to independently learn representations for the normalized scattering and incandescent channels.



3 Analysis of Learned Latent Representations of the L-II signals

We first use the VAE to encode the L-II signals from the training data set for both Ch.0 and Ch.1 into lower dimensional latent representations, which we refer to as z_i . We refer to the latent variables for Ch. 0 as z_1 and z_2 , and the latent variables for Ch. 1 as z_3 and z_4 . Because these latent variables are learned representations of the normalized Ch.0 and Ch.1 time series, these variables provide information about the shape of these signals. Thus, the distance between variables in the latent space representations provide a means to visualize how similar L-II signals are to one another.

The smoothness constraint in the VAE ensures that signals with greater similarity are mapped to points that are closer together in the latent space, preserving meaningful structure in the learned representation. This constraint encourages continuity, meaning that small changes in the input signal result in gradual variations in the latent representation. However, the specific structure of these learned representations is inherently shaped by the underlying data distribution, as the VAE optimizes its encoding based on the patterns present in the underlying data. The distribution of latent variables for each aerosol class can provide information about how similar signals within each class are to one another, and it can also provide information about how much separability there is between different aerosol classes in terms of their latent space representations. These distributions provide insights into whether the shape of the normalized signals for Ch. 0 and Ch. 1 can be used to differentiate classes of absorbing aerosols that the SP2 is sensitive to.

The latent variables learned from the normalized Ch. 0 and Ch. 1 signals demonstrate smoothly varying characteristics when we plot the distributions of the aerosol populations in terms of I_{max} vs. CR (Figure 3). I_{max} (proportional to the mass of the refractory portion of the aerosol) and CR (proportional to the temperature at which an aerosol incandesces in the SP2's laser beam) are strongly correlated with both the size of the aerosol and its chemical composition. As discussed in Section 2.2, the left mode in Figure 3 is typical of FeO_x aerosols, and the right mode is typical of rBC aerosols and their proxies, including FS, due to differences in their characteristic incandescence temperatures. Since the latent representations are smoothly varying in this space, this suggests that the shapes of both the scattering and incandescent signals are strongly correlated with their overall distributions in terms of refractory aerosol mass and incandescence temperature. This suggests that the shape of the L-II signals alone (without information about the signal magnitude) provides information about the physio-chemical properties of the aerosols detected in the SP2.

To better understand how the latent variables represent the shapes of the L-II signals, we use the trained decoders from each VAE to map out the latent space representations in terms of representative Ch. 0 and Ch. 1 time series. Figure 4 shows the latent manifold for Ch. 0, left, and for Ch. 1, right. By generating characteristic signals along the deciles of the distributions of the latent variables for Ch. 0 and Ch. 1, we can examine how the latent variables capture specific properties of the time series signals such as symmetry and saturation. The saturation of signals in the Ch. 0 time series (Figure 4, left panel) shows flatness at the peaks, which is most common for dust-like particles (ATD, VA, and CFA), as the scattering detector becomes saturated for these large particles, and the peaks are artificially flattened. These signals are more evident in the bottom right side of the latent manifold for Ch. 0. The iron oxide aerosols Fe_2O_3 and Fe_3O_4 commonly have symmetric scattering signals,

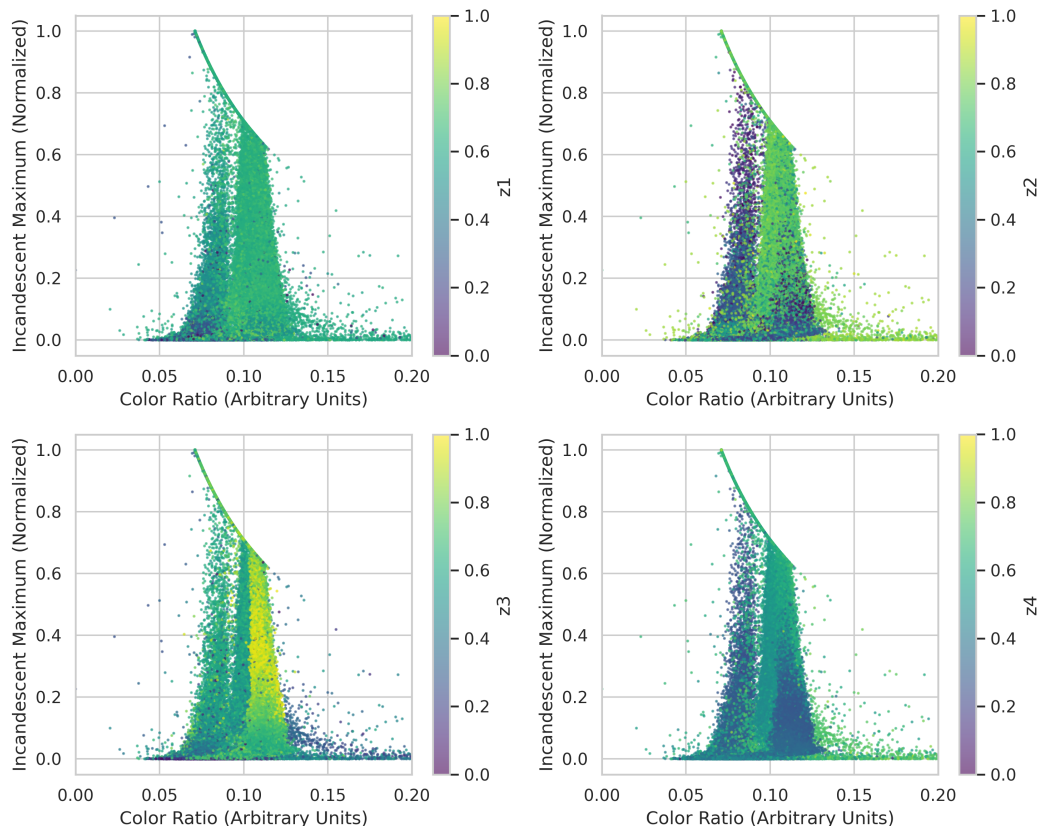


Figure 3. Incandescent-peak-height to color-ratio relationship for different incandescent aerosols visualized by latent representation: Each point represents a single aerosol particle, and we color the points by the latent space representations for Ch. 0 (z_1 and z_2 , top panels) and Ch. 1 (z_3 and z_4 , bottom panels). Here, we normalize z_1 , z_2 , z_3 and z_4 between 0 and 1, since only the relative value of the learned latent variables is meaningful.

230 which is evident in the upper right of the latent manifold for Ch. 0. The latent manifold for Ch. 1 captures the symmetry of the incandescent signal along one dimension, and the narrowness of the incandescent peak along the other dimension.

Since the size of the latent representations learned by the VAE is a hyper-parameter, we also explore how varying the size of the latent space impacts our ability to differentiate aerosols by class using these lower dimensional latent representations. To do this, we train VAE's with a latent vector z that has either $n=2$ and $n=3$ variables. We then visualize the distributions of the latent variables from our training data set for each of the 7 aerosol classes represented in our dataset. Figure 5 shows the distribution of the encodings of Ch.0 (top) and Ch. 1 (bottom) when training two VAE's with $n=2$ variables in its latent space representation, and Figure 6 shows the distributions of the encodings of Ch. 0 (top) and Ch. 1 (bottom) when training two VAE's with $n=3$ variables. Each point on these plots represents the L-II signal from an aerosol detected in the SP2, encoded by

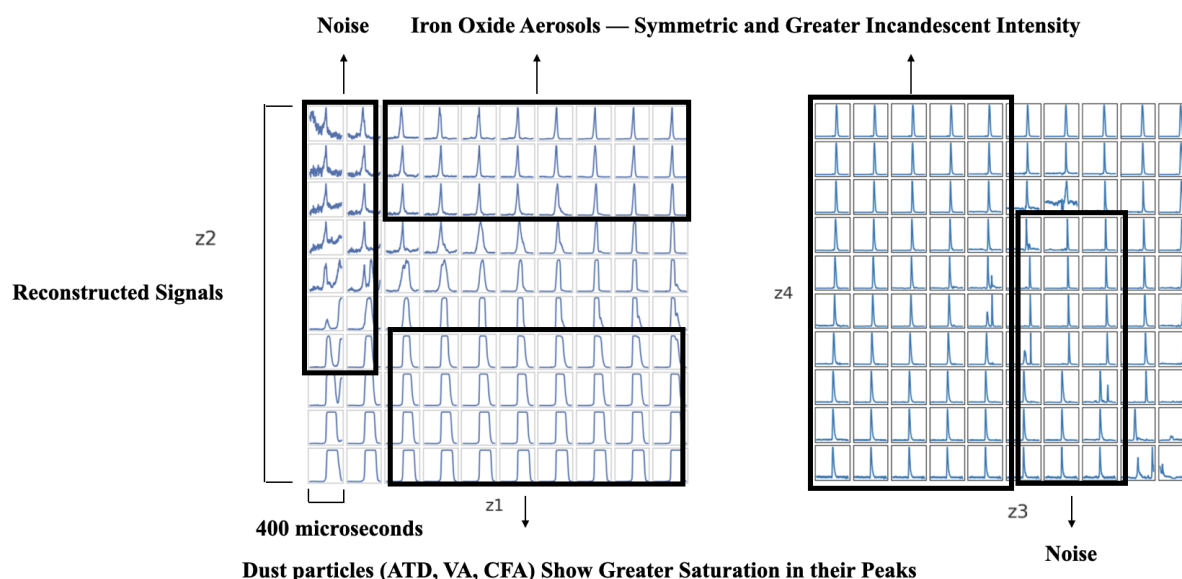


Figure 4. Latent Manifolds for the L-II Signals. Left: Latent Manifold for Channel 0. Right: Latent Manifold for Channel 1.

the VAE into its latent space representation. These latent space representations therefore gives us a useful way to visualize the
240 distributions of L-II signals found in each class.

In examining the latent space distributions for $n=2$ (Figure 5), we can find some consistency in terms of the encodings for
different aerosol classes. For the scattering channel (Ch. 0) and incandescent channel (Ch. 1), the latent space representations
for the black carbon proxies (FS and FS+glyc) show significant overlap, as do the latent space encodings for the FeO_x proxies
(Fe_2O_3 and Fe_3O_4). Regions of high density in these latent space representations indicate more examples of that L-II signal
245 are found when measuring that aerosol class. The latent space representations for the dust-like aerosols (CFA, ATD, and VA)
show less clear regions of high density when compared to the black carbon proxies and FeO_x proxies, which is explained by
the greater variability across classes in observed L-II signals for the dust-like aerosols. In Lamb (2019), it was noted that these
dust-like aerosols were more likely to lead to saturated signals (due to their large optical size) or other irregularities in their
L-II signals when compared with rBC or FeO_x . This greater variability in L-II signals is reflected in the greater spread in their
250 latent space representations.

Similarly, the latent space distributions for $n=3$ (Figure 6) show some consistency across aerosols of similar classes. How-
ever, the additional dimension in the latent space representations provides additional contrast between aerosols of similar
classes. For example, for FS+glyc, the latent space representation for Ch.1 in part significantly overlaps with that for uncoated
FS, but also shows a large density of points in space that is not represented at all by the FS signals. In addition, while VA,
255 ATD, and CFA show similarities in their latent space representations for Ch. 0, there are more clear differences in the density
of points observed in their latent space representations for Ch. 1, suggesting that the incandescent signal provides greater con-

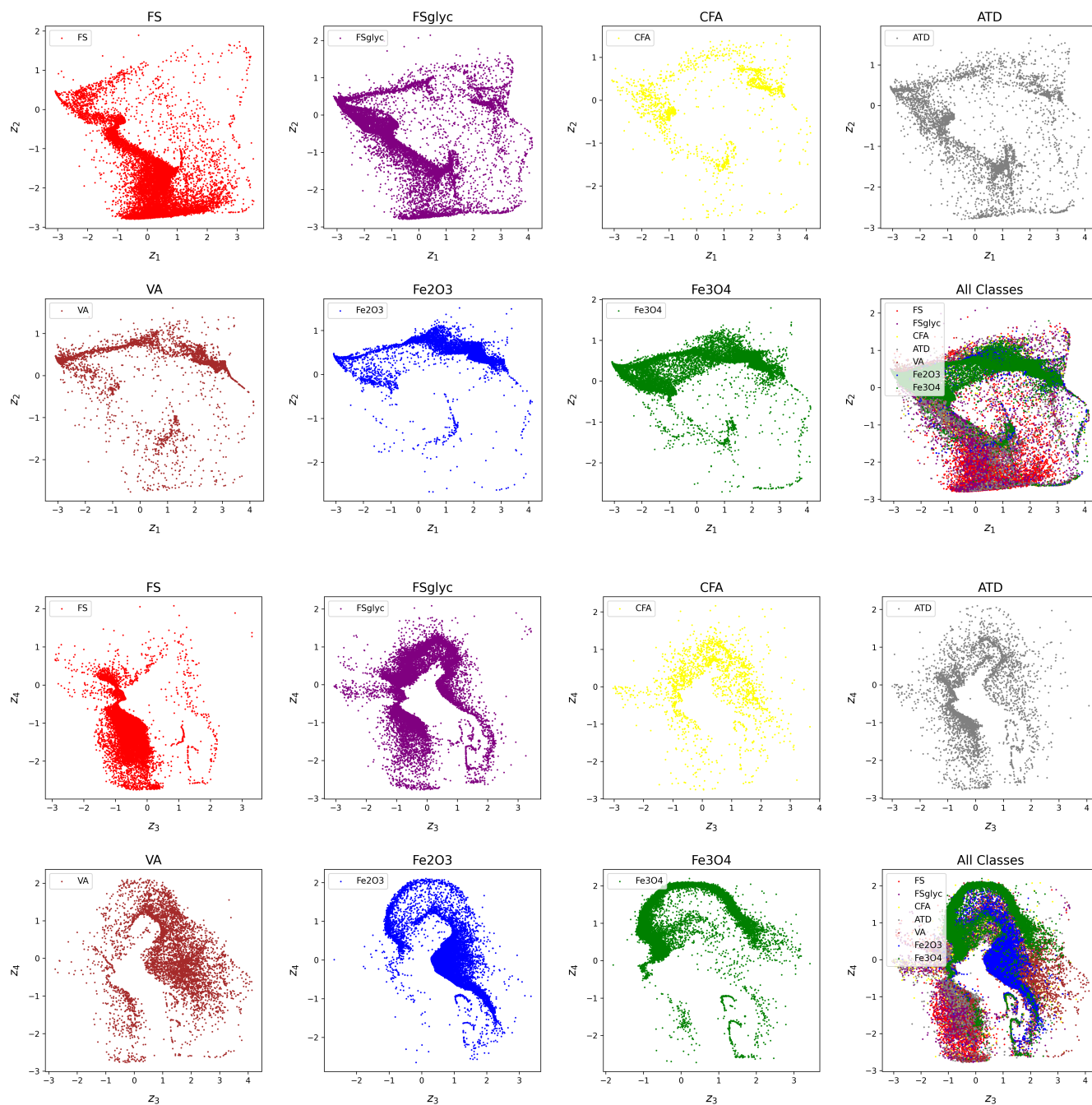


Figure 5. Latent Space for $n = 2$. Above: Latent space distribution for Ch. 0 (scattering channel) shown by each aerosol class. Below: Latent space distribution for Ch 1 (incandescent channel) shown by each aerosol class.

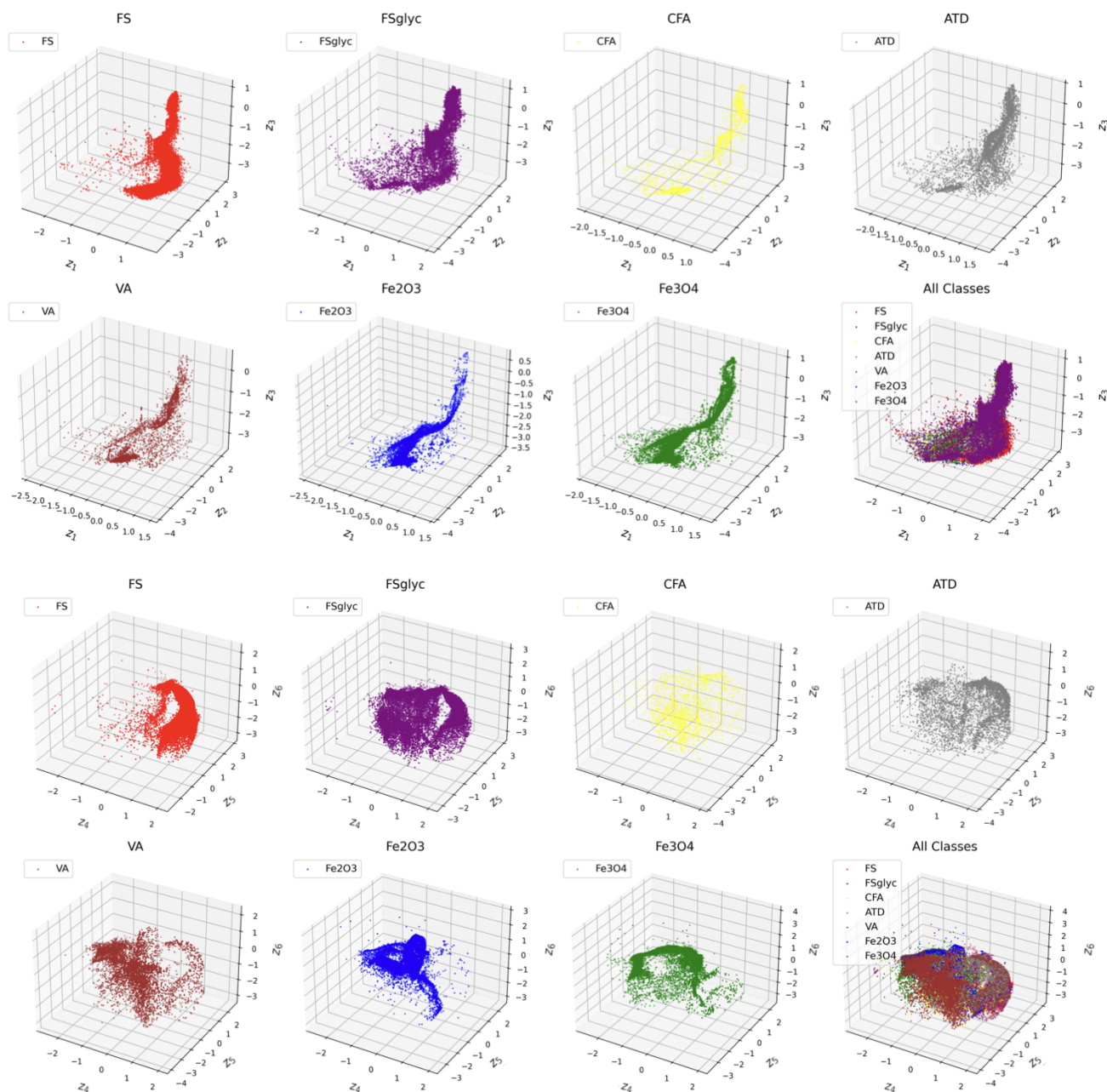


Figure 6. Latent Space for $n = 3$. Above: Latent space distribution for Ch. 0 (scattering channel) shown by each aerosol class. Below: Latent space distribution for Ch 1 (incandescent channel) shown by each aerosol class.



trast in differentiating the signals of these dust-like aerosols by class. Physically, this makes sense, since the L-II incandescent signal could provide greater contrast between metallic inclusions of different chemical composition or size than the scattering signal, and these likely differ between the ATD, VA, and CFA populations. Finally, the iron oxide aerosols also show more clear separability in terms of their latent space representations for Ch. 1.

We found that a latent space with $n=2$ provided sufficient separability and interpretability of latent space variables of the L-II time series, compared to higher dimensional latent space representations. Increasing the size of the latent space enables the VAE to capture a higher level of detail and complexity in the original L-II signals, allowing for more accurate and nuanced reconstructions. However, high dimensionality for the latent space representation is more difficult to interpret, and may not provide additional meaningful information for downstream tasks. We explored a latent space with $n=4$ and $n=5$ (not shown), but found that including additional variables did not provide improved performance on downstream tasks like classification (Section 5), suggesting that the majority of the variance in the L-II signals can be encoded into 2 latent dimensions.

4 Identifying Outliers in Aerosol Populations using the L-II Latent Space Representations

In addition to visualizing the variability of the populations of aerosols detected by the SP2, the latent space representation can also be used as a method to identify outliers in populations of aerosols. To illustrate this, we define a metric to identify outliers based on their distance from the global centroid of the latent space distribution.

First, we calculate the global centroid of the latent space representations that have been encoded into the latent space. Using the Euclidean distance, we then determine the mean straight-line distance between the latent space representations for each of the L-II signals for individual aerosols and the center of the latent space manifold, μ_c , as well as the standard deviation of this distance, σ_c . Outliers can be defined by setting a sensitivity threshold, ϵ , such that latent space representations that are further away from the global centroid than σ_c multiplied by this threshold are considered outliers:

$$d_{outlier} = \mu_c + \epsilon \times \sigma_c \quad (7)$$

As the sensitivity threshold ϵ increases, the number of L-II signals identified as outliers increases. With this metric, we can identify outliers in our latent space representation as shown in Figure 7 for the 2D latent space representation for Ch. 1. We find that the latent space representations that lie further from the center of the latent space manifold are more likely to exhibit spurious behavior in the L-II signals. We show two examples of L-II signals that were identified as outliers and their latent space locations; one of these L-II signals is from an Fe_3O_4 particle (top panel, right), and the other is CFA (bottom panel, right). From looking at the original L-II signals, it is clear why these two aerosols were identified as outliers. Both shows scattering peaks that are significantly off-center from the typical triggering location (Ch. 0), and they also demonstrate significantly later incandescence (Ch. 1).

One challenge with using the SP2 to identify FeO_x aerosols is that some dust-like aerosols have similar I_{max} vs. CR values as FeO_x (Lamb, 2019; Lamb et al., 2021). This makes it challenging to quantify FeO_x mass loadings in atmospheric conditions, particularly in remote regions where particles need to be differentiated on an individual, rather than population, basis (Lamb

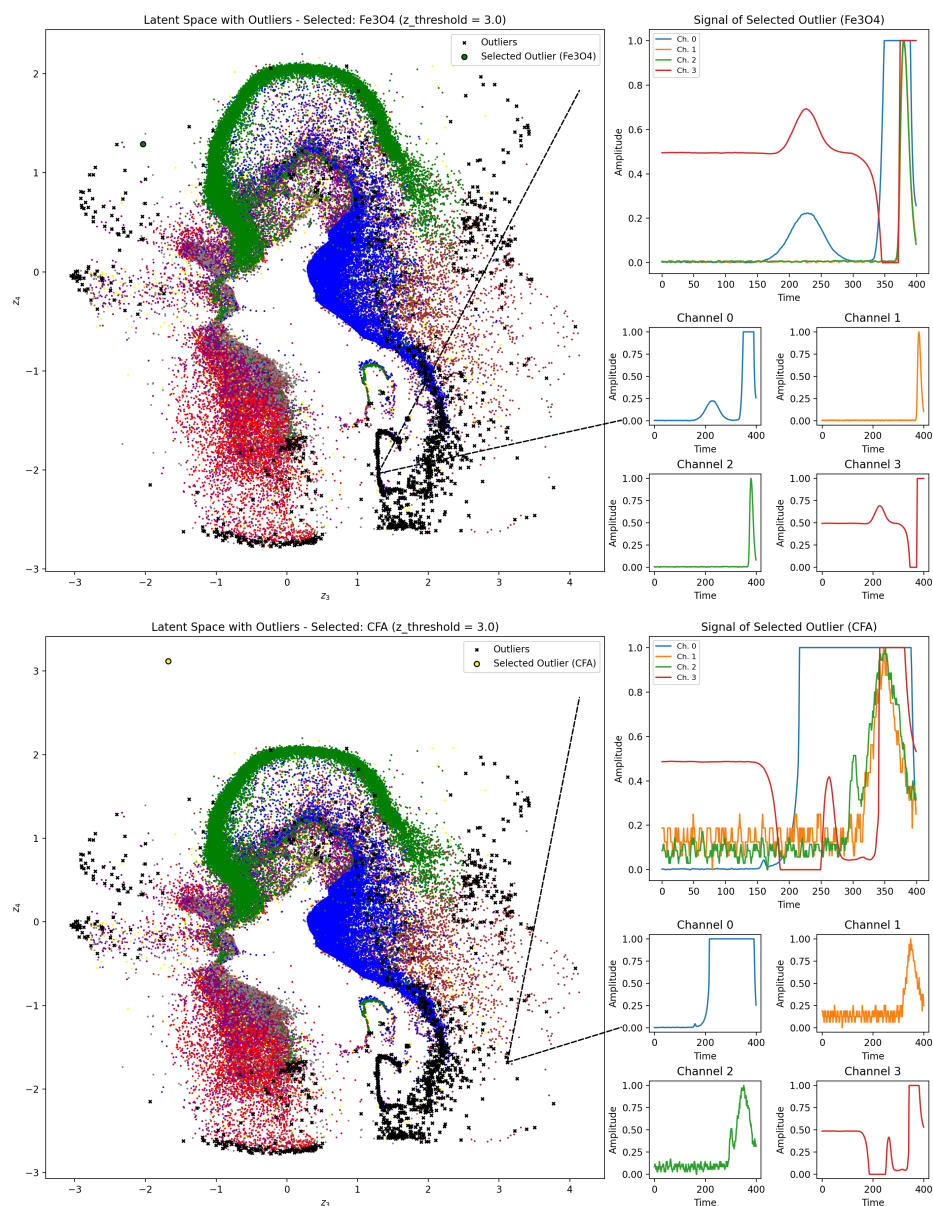


Figure 7. Examples of outlier detection via latent space representation. Black points indicate points determined to be outliers based on their Euclidean distance from the center of the 2D Latent Space for Channel 1. Examples of the L-II signals for an Fe₃O₄ aerosol (top, right) and a CFA aerosol (bottom, right) that have been identified as outliers based on the distance of their latent space embedding from the center of the latent space are shown.

et al., 2021). The outlier detection method using the latent space representations could be used as an alternative approach to
 290 identify dust-like aerosols that may be mis-categorized as FeO_x. Due to the random, scattered nature of the L-II signals for



dust-like particles, such a threshold can minimize the occurrence of false positives when identifying ambient aerosols as FeO_x , and thus improve the accuracy of FeO_x mass loading measurements.

Furthermore, outlier detection could be used as an approach to identify interesting populations of aerosols in ground-based or airborne field observational datasets in an unsupervised manner. A number of recent studies have used observations from the SP2 to identify unique populations of aerosols from their L-II signals, including tar brown carbon (Corbin and Gysel-Beer, 2019), iron oxide aerosols (Lamb et al., 2021), or black carbon associated with pyrocumulonimbus (Katich et al., 2023). By automatizing the detection of interesting or unique populations of aerosols that the SP2 detects, this method can help to identify or characterize aerosols from different atmospheric sources that may not be evident with traditional SP2 L-II analysis approaches. The latent space representations may also provide additional information about the shape of the L-II signals (beyond coating state or refractory aerosol mass) that can be linked to the physio-chemical characteristics of the aerosol particles.

5 Classifying aerosols using their L-II latent space representations

As a final down-stream task, we investigate how useful the latent space representations of the L-II time series are in terms of differentiating absorbing aerosols that the SP2 is sensitive to by class. To do this, we follow an approach similar to the supervised classification approach, using a random forest algorithm to classify aerosols in a supervised manner, as described in Lamb (2019). However, rather than doing significant feature engineering on the L-II signals as input features to train the random forest, we instead use the learned latent space representations from Ch. 0 and Ch. 1.

For each sampled aerosol, we concatenate the latent variables from Ch. 0 and Ch. 1, as well as the S_{max} , I_{max} , and CR derived from the L-II signal for that aerosol (Figure 2). We test using both 2 latent variables each from Ch. 0 and Ch.1 and 3 latent variables. In the follow text, we refer to the case with the 2 latent variables each from Ch. 0 and Ch. 1 as the “2D latent space RF” and the 3 latent variables each from Ch. 0 and Ch. 1 as the “3D latent space RF”.

The Random Forest algorithm constructs an ensemble of decision trees, using a subset of the training samples to construct each decision tree. The class for each sample is then determined by the vote of all the randomly constructed decision trees. Here, we use 67% of the data for training and 33% for testing the random forest, and use the Random Forest algorithm as implemented in the scikit-learn package (Pedregosa et al., 2011).

The confusion matrices for the results of the random forest classification for the 2D latent space RF (left) and the 3D latent space RF (right) is shown in Figure 8. The 3D latent space RF and 2D latent space RF demonstrated similar performance. We also investigated 4D and 5D latent space representations as input to the RF, and did not find significant further improvements over the 2D or 3D cases, suggesting that 2 or 3 latent variables already provides the majority of meaningful information from the L-II signals in terms of differentiating absorbing aerosols.

Compared with the random forest that used significant feature engineering (Lamb, 2019), the 2D and 3D latent space RF’s are able to more accurately classify aerosols of each of the 3 major sub-types that the SP2 detects (BC, dust-like, and FeO_x). In particular, we find that the 2D and 3D latent space RF’s are able to do a significantly better job of differentiating Fe_2O_3 and



Fe₃O₄ from one another; in Lamb (2019) the RF with significant feature engineering mis-identified Fe₂O₃ as Fe₃O₄ 50% of the time. This is likely due to the differences that are evident in the latent space representations from Ch. 1 (Figure 5), indicating that there are differences in the incandescent signals for Fe₂O₃ and Fe₃O₄ that are not readily evident from a simple analysis of the CR. This suggests that there are meaningful differences between iron oxide aerosols in terms of their L-II signals, which could be further exploited to improve the detection accuracy of these aerosols using the SP2, and warrants further investigation. This result further illustrates how the learned latent space representations can identify new information from the L-II signals that is not obvious following traditional SP2 analysis approaches. The 2D and 3D latent space RFs also do a slightly better job at differentiating the 3 different classes of dust-like aerosols from one another, particularly for VA, compared with the RF with significant feature engineering described in Lamb (2019).

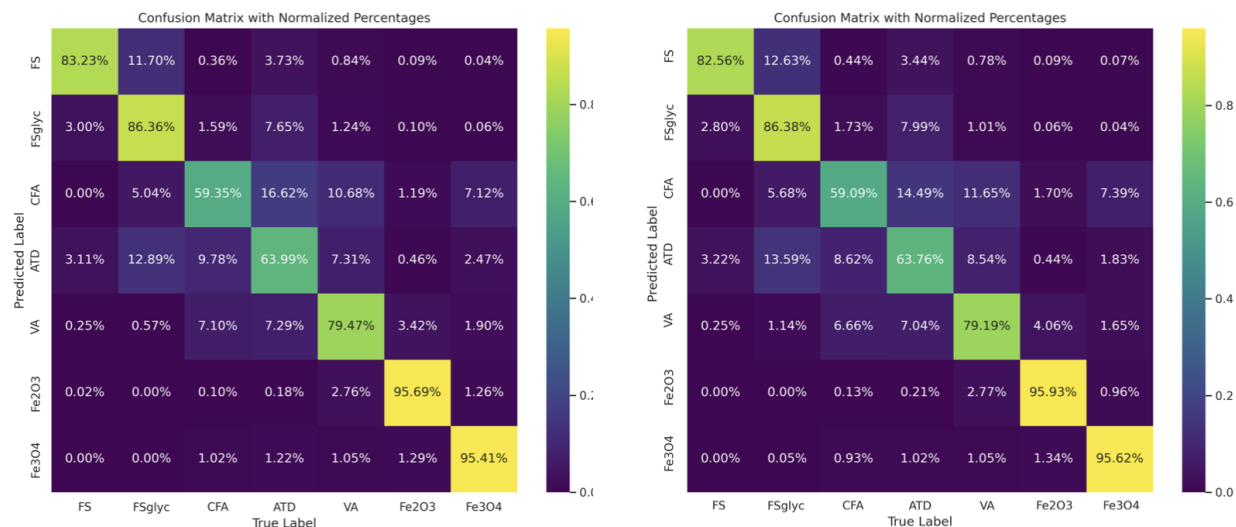


Figure 8. Performance of RF on classifying different types of absorbing aerosols observed by the SP2 based on their latent representations. Left: Confusion matrix for 2D latent space RF. Right: Confusion matrix for the 3D latent space RF.

6 Conclusions

In this paper, we have built on our previous research in Lamb (2019) to further explore how data-driven methods such as machine learning can be applied to L-II time series from the SP2 instrument. Here we have focused on unsupervised machine learning as a method to classify aerosols and discover more of their chemio-physical properties, thus demonstrating a path towards better understanding the variability of aerosols observed in the atmosphere. By using a variational autoencoder to encode the L-II signals from the SP2, we found that this approach could be used to quantify the variability within aerosol classes. It can also be used to identify outliers in observational data sets, and is a promising method for identifying unique or interesting aerosol populations automatically in large data-sets from research field campaigns. In addition, we were able to



achieve high separability between the aerosols classes when using the latent space representations of the L-II signals as input to a supervised classification algorithm. In particular the distinct separation observed between the two classes of iron oxide aerosols was promising, suggesting that differences between the response of these two iron oxide aerosols in the SP2 could be further exploited to improve the detectability of these aerosols.

345 Black carbon and other aerosols that are detected by the SP2, are operationally defined– that is, we classify aerosol by composition based on their response in the instrument and their similarity to laboratory-based proxies for these aerosols. However, populations of aerosols in the atmosphere, even those from the same emission sources, will have a distribution of characteristics in terms of their composition, optical properties, and sizes. These differences contribute to variations in the response of these aerosols in the SP2. The unsupervised machine learning approach that we discuss in this paper is a method
350 where these types of variations between L-II signals can be quantified and assessed in a lower dimensional representation space.

There are a number of promising future research directions in terms of applying unsupervised and semi-supervised machine learning methods to the L-II signals from the SP2. In particular, contrastive learning is a promising approach for identifying latent space representations that most meaningfully separate aerosols of different classes (Severson et al., 2019; Abid and Zou,
355 2019). This approach could improve classification when training on labeled laboratory data sets and applying to atmospheric data sets. Similarity metrics, such as Euclidean distance, can be used to determine how similar the SP2 response of aerosols are, enabling more quantitative comparison across field campaign data sets and regions of the atmosphere (Levy et al., 2024). Future research should also further explore how instrument configuration will impact learned representations from the SP2 signal, such that observational data sets across field campaigns and from different instruments could be meaningfully combined towards
360 improved SP2 analysis.

Code and data availability. Code to reproduce the analysis in this paper is available at <https://github.com/adoshi25/SP2-Aerosol-Classification>. Labeled, machine learning ready data sets for the SP2 L-II signals are available at <https://doi.org/10.5281/zenodo.15800436>.

Author contributions. Conceptualization: K.D.L. Methodology: A.D; K.D.L. Data curation: K.D.L. Data visualization: A.D. Writing original draft: A.D; K.D.L. All authors approved the final submitted draft.

365 *Competing interests.* The authors declare they have no competing interests.

Acknowledgements. K.D.L. acknowledges support from the Zegar Family Foundation and the NSF LEAP Center at Columbia University.



References

- Abid, A. and Zou, J.: Contrastive variational autoencoder enhances salient features, arXiv preprint arXiv:1902.04601, 2019.
- Adachi, K., Moteki, N., Kondo, Y., and Igarashi, Y.: Mixing states of light-absorbing particles measured using a transmission electron
370 microscope and a single-particle soot photometer in Tokyo, Japan, *Journal of Geophysical Research: Atmospheres*, 121, 9153–9164, 2016.
- Baumgardner, D., Popovicheva, O., Allan, J., Bernardoni, V., Cao, J., Cavalli, F., Cozic, J., Diapouli, E., Eleftheriadis, K., Genberg, P., et al.: Soot reference materials for instrument calibration and intercomparisons: a workshop summary with recommendations, *Atmospheric Measurement Techniques*, 5, 1869–1887, 2012.
- 375 Bingham, E., Chen, J. P., Jankowiak, M., Obermeyer, F., Pradhan, N., Karaletsos, T., Singh, R., Szerlip, P., Horsfall, P., and Goodman, N. D.: Pyro: Deep Universal Probabilistic Programming, <https://doi.org/10.48550/arXiv.1810.09538>, arXiv:1810.09538 [cs], 2018.
- Biswas, A., Ziatdinov, M., and Kalinin, S. V.: Combining variational autoencoders and physical bias for improved microscopy data analysis*, *Machine Learning: Science and Technology*, 4, 045 004, <https://doi.org/10.1088/2632-2153/acf6a9>, publisher: IOP Publishing, 2023.
- Bond, T. C., Doherty, S. J., Fahey, D. W., Forster, P. M., Berntsen, T., DeAngelo, B. J., Flanner, M. G., Ghan, S., Kärcher, B., Koch, D.,
380 et al.: Bounding the role of black carbon in the climate system: A scientific assessment, *Journal of geophysical research: Atmospheres*, 118, 5380–5552, 2013.
- Corbin, J. C. and Gysel-Beer, M.: Detection of tar brown carbon with a single particle soot photometer (SP2), *Atmospheric Chemistry and Physics*, 19, 15 673–15 690, 2019.
- Gao, R., Schwarz, J., Kelly, K., Fahey, D., Watts, L., Thompson, T., Spackman, J., Slowik, J., Cross, E., Han, J.-H., et al.: A novel method for
385 estimating light-scattering properties of soot aerosols using a modified single-particle soot photometer, *Aerosol Science and Technology*, 41, 125–135, 2007.
- Heimerl, K., Weinzierl, B., Gysel, M., Baumgardner, D., Kok, G., Linke, C., Schnaiter, M., Schwarz, J., Sheridan, P., Subramanian, R., et al.: Using a Single Particle Soot Photometer to detect and distinguish different absorbing aerosol types, in: *European Aerosol Conference*, vol. 44, pp. 663–675, European Aerosol Conference, 2012.
- 390 Katich, J., Apel, E., Bourgeois, I., Brock, C., Bui, T., Campuzano-Jost, P., Commane, R., Daube, B., Dollner, M., Fromm, M., et al.: Pyrocumulonimbus affect average stratospheric aerosol composition, *Science*, 379, 815–820, 2023.
- Kingma, D. P. and Ba, J.: Adam: A Method for Stochastic Optimization, <https://doi.org/10.48550/arXiv.1412.6980>, arXiv:1412.6980 [cs], 2017.
- Kingma, D. P. and Welling, M.: Auto-Encoding Variational Bayes, <https://doi.org/10.48550/arXiv.1312.6114>, arXiv:1312.6114 [stat], 2022.
- 395 Lamb, K.: Laser-Induced Incandescent Signals for Laboratory Samples of Absorbing Aerosols Detected by the Single Particle Soot Photometer, [Data set]. Zenodo., <https://doi.org/10.5281/zenodo.15800436>, 2025.
- Lamb, K., Matsui, H., Katich, J., Perring, A., Spackman, J., Weinzierl, B., Dollner, M., and Schwarz, J.: Global-scale constraints on light-absorbing anthropogenic iron oxide aerosols, *Npj Climate and Atmospheric Science*, 4, 15, 2021.
- Lamb, K. D.: Classification of iron oxide aerosols by a single particle soot photometer using supervised machine learning, *Atmospheric*
400 *Measurement Techniques*, 12, 3885–3906, 2019.
- Lamb, K. D., Perring, A. E., Samset, B., Peterson, D., Davis, S., Anderson, B. E., Beyersdorf, A., Blake, D. R., Campuzano-Jost, P., Corr, C. A., et al.: Estimating source region influences on black carbon abundance, microphysics, and radiative effect observed over South Korea, *Journal of Geophysical Research: Atmospheres*, 123, 13–527, 2018.



- Levy, A., Shalom, B. R., and Chalamish, M.: A Guide to Similarity Measures, arXiv preprint arXiv:2408.07706, 2024.
- 405 Liu, D., Taylor, J. W., Crosier, J., Marsden, N., Bower, K. N., Lloyd, G., Ryder, C. L., Brooke, J. K., Cotton, R., Marengo, F., et al.: Aircraft and ground measurements of dust aerosols over the west African coast in summer 2015 during ICE-D and AER-D, *Atmospheric Chemistry and Physics*, 18, 3817–3838, 2018.
- Moteki, N. and Kondo, Y.: Dependence of laser-induced incandescence on physical properties of black carbon aerosols: Measurements and theoretical interpretation, *Aerosol Science and Technology*, 44, 663–675, 2010.
- 410 Moteki, N., Adachi, K., Ohata, S., Yoshida, A., Harigaya, T., Koike, M., and Kondo, Y.: Anthropogenic iron oxide aerosols enhance atmospheric heating, *Nature Communications*, 8, 15 329, 2017.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al.: Scikit-learn: Machine learning in Python, the *Journal of machine Learning research*, 12, 2825–2830, 2011.
- Schwarz, J., Gao, R., Fahey, D., Thomson, D., Watts, L., Wilson, J., Reeves, J., Darbeheshti, M., Baumgardner, D., Kok, G., et al.: Single-
415 particle measurements of midlatitude black carbon and light-scattering aerosols from the boundary layer to the lower stratosphere, *Journal of Geophysical Research: Atmospheres*, 111, 2006.
- Schwarz, J., Spackman, J., Gao, R., Perrington, A., Cross, E., Onasch, T., Ahern, A., Wrobel, W., Davidovits, P., Olfert, J., et al.: The detection efficiency of the single particle soot photometer, *Aerosol Science and Technology*, 44, 612–628, 2010.
- Severson, K. A., Ghosh, S., and Ng, K.: Unsupervised learning with contrastive latent variable models, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 4862–4869, 2019.
- 420 Stephens, M., Turner, N., and Sandberg, J.: Particle identification by laser-induced incandescence in a solid-state laser cavity, *Applied optics*, 42, 3726–3736, 2003.
- Wei, R., Garcia, C., El-Sayed, A., Peterson, V., and Mahmood, A.: Variations in Variational Autoencoders - A Comparative Evaluation, *IEEE Access*, 8, 153 651–153 670, <https://doi.org/10.1109/ACCESS.2020.3018151>, conference Name: IEEE Access, 2020.
- 425 Yoshida, A., Moteki, N., Ohata, S., Mori, T., Tada, R., Dagsson-Waldhauserová, P., and Kondo, Y.: Detection of light-absorbing iron oxide particles using a modified single-particle soot photometer, *Aerosol Science and Technology*, 50, 1–4, 2016.
- Ziatdinov, M.: PyroVED, <https://github.com/ziatdinovmax/pyroVED>.