
LATENT DIRICHLET TRANSFORMER VAE FOR HYPERSPECTRAL UNMIXING WITH BUNDLED ENDMEMBERS

Giancarlo Giannetti

Ontario Tech University, Canada
giancarlo.giannetti@ontariotechu.net

Faisal Z. Qureshi

Ontario Tech University, Canada
faisal.qureshi@ontariotechu.ca

ABSTRACT

Hyperspectral images capture rich spectral information that enables per-pixel material identification; however, spectral mixing often obscures pure material signatures. To address this challenge, we propose the Latent Dirichlet Transformer Variational Autoencoder (LDVAE-T) for hyperspectral unmixing. Our model combines the global context modeling capabilities of transformer architectures with physically meaningful constraints imposed by a Dirichlet prior in the latent space. This prior naturally enforces the sum-to-one and non-negativity conditions essential for abundance estimation, thereby improving the quality of predicted mixing ratios. A key contribution of LDVAE-T is its treatment of materials as bundled endmembers, rather than relying on fixed ground truth spectra. In the proposed method our decoder predicts, for each endmember and each patch, a mean spectrum together with a structured (segmentwise) covariance that captures correlated spectral variability. Reconstructions are formed by mixing these learned bundles with Dirichlet-distributed abundances garnered from a transformer encoder, allowing the model to represent intrinsic material variability while preserving physical interpretability. We evaluate our approach on three benchmark datasets, Samson, Jasper Ridge, and HYDICE Urban and show that LDVAE-T consistently outperforms state-of-the-art models in abundance estimation and endmember extraction, as measured by root mean squared error and spectral angle distance, respectively.

1 Introduction

Hyperspectral images capture data across a wide range of wavelengths, typically over 100, spanning the infrared, visible, and sometimes ultraviolet spectra. Each material within a scene exhibits a unique spectral signature, enabling per-pixel material identification through appropriate data processing. In satellite-based remote sensing, however, each pixel often corresponds to a large ground area containing multiple materials. This results in spectral mixing, which necessitates hyperspectral unmixing, the process of decomposing mixed pixel spectra into their constituent endmembers.

Hyperspectral Unmixing (HU) involves decomposing a hyperspectral image into its constituent materials, referred to as endmembers, and estimating the relative abundance of each material within every pixel. This task is essential in many remote sensing applications, as each pixel often represents a mixture of materials due to limited spatial resolution. Traditional HU methods are broadly categorized into linear and nonlinear models. Linear mixing models (LMMs) assume that the observed spectrum is a convex combination of endmember signatures, weighted by their abundance fractions. While computationally efficient and physically interpretable, LMMs often fail to capture complex light interactions such as scattering and nonlinearity, especially in densely vegetated or urban areas.

Nonlinear unmixing models aim to address these limitations by incorporating physical or data-driven nonlinearities into the mixing process. However, these approaches typically involve greater model complexity and higher computational cost. More recently, machine learning and deep learning methods have emerged as powerful tools for HU, offering improved flexibility in modeling spectral and spatial patterns. These include autoencoders, variational models, graph-based networks, and transformer architectures, many of which incorporate priors or constraints (e.g.,

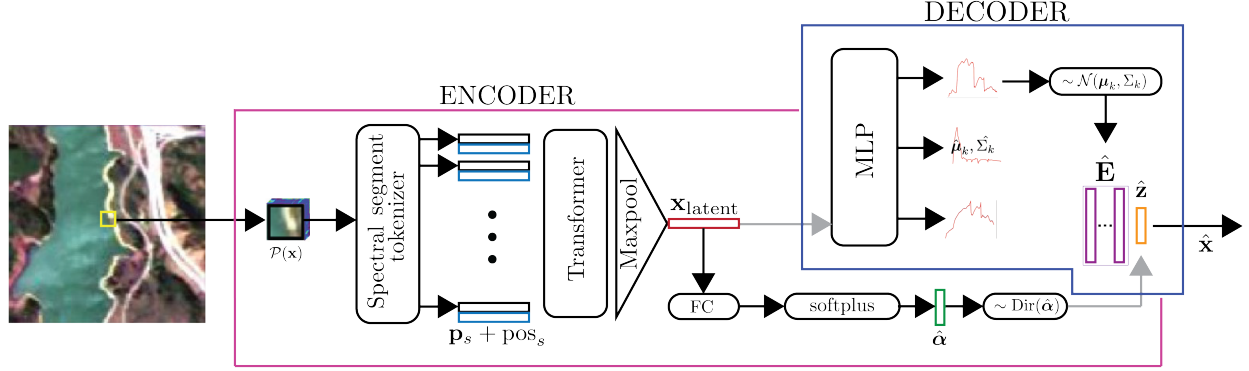


Figure 1: Overview. Transformer based encoder takes a hyperspectral image patch $\mathcal{P}(\mathbf{x})$ and construct $\mathbf{x}_{\text{latent}}$, $\hat{\alpha}$ and abundances $\hat{\mathbf{z}}$. The decoder takes $\mathbf{x}_{\text{latent}}$ and $\hat{\mathbf{z}}$ and reconstructs pixel \mathbf{x} . The decoder also constructs endmembers $\hat{\mathbf{e}}_k$.

sparsity, non-negativity, sum-to-one) to enhance physical interpretability. Despite these advances, challenges remain in generalizing across diverse scenes, extracting pure endmembers from limited ground truth, and balancing model expressiveness with physical realism. As such, HU continues to be an active area of research at the intersection of remote sensing, signal processing, and machine learning.

State-of-the-art (SOTA) methods increasingly rely on Vision Transformer (ViT) architectures [1–4], which outperform traditional multilayer perceptron (MLP) and convolutional neural network (CNN) approaches by effectively capturing long-range dependencies [1], a critical capability for estimating abundances in hyperspectral data. In contrast, CNNs are inherently limited to local neighborhood information, which may be insufficient for modeling complex spectral-spatial interactions.

To better address the unique characteristics of hyperspectral images, such as locally homogeneous regions, nonlinear mixing effects, and limited labeled data—researchers have proposed several specialized modules to enhance unmixing performance. One promising direction involves incorporating Dirichlet distributions into a model’s latent space [5, 6], enabling Variational Autoencoders (VAEs) to better model the constraints of hyperspectral unmixing. However, such Dirichlet-based approaches have thus far been limited to MLP and CNN backbones, and have not yet been explored in transformer-based architectures. To bridge this gap, we propose LDVAE-T, which integrates Dirichlet priors into a ViT-based VAE framework, combining the probabilistic interpretability of Dirichlet modeling with the global feature extraction capabilities of transformers. Our key novelty is a “bundled endmember” decoder: instead of using fixed spectra, the model predicts, for each endmember and for each patch, a distributional prototype composed of a mean spectrum and a structured (segmentwise) covariance that captures correlated spectral variability. Reconstructions are then formed by mixing samples from these learned bundles according to Dirichlet-distributed abundances. Further, we tokenize spectra as fixed-length segments to respect spectral locality while enabling efficient attention. Together, these design choices allow LDVAE-T to model intrinsic material variability while preserving the sum-to-one and non-negativity properties central to HU, leading to improved endmember extraction and abundance estimation across diverse scenes.

2 Related Work

Pixel unmixing approaches in hyperspectral imaging can be broadly categorized into two groups: (a) physics-based methods and (b) data-driven techniques. Physics-based methods rely on models that describe how light interacts with materials, such as Hapke’s Bidirectional Reflectance Distribution Function (BRDF)[7] and the Atmospheric Dispersion Model [8]. While these approaches offer physically grounded interpretations, they often require detailed, scene-specific radiative parameters, limiting their practicality in real-world applications. In contrast, data-driven methods dominate the field due to their flexibility and ease of deployment, though they are highly dependent on the availability of annotated training data. Hybrid models that fuse physics-based insights with data-driven learning have also been studied [9]. A widely adopted class of data-driven methods is based on Non-Negative Matrix Factorization (NMF), which models the hyperspectral image as the product of two matrices: one representing endmembers and the other representing abundances. Several extensions to NMF based methods have been proposed in the literature that exploit spatial or spectral neighbourhood structure [10–12], use iterative refinement to improve results [13], or integrate handcrafted and learned priors for improved generalization [14].

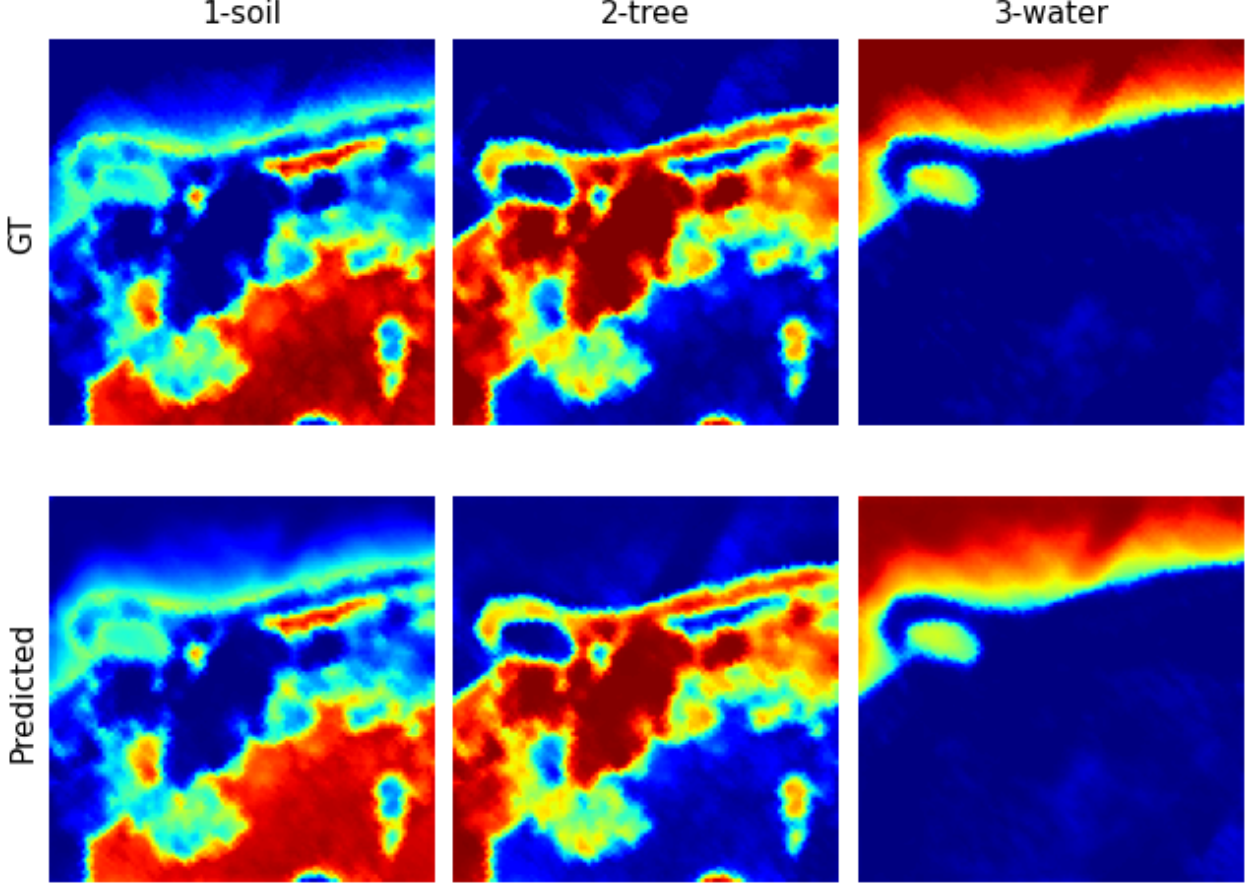


Figure 2: Ground Truth vs Predicted endmember heat-maps for Samson dataset. These plots visualize per-pixel abundances.

Recent efforts have shifted toward deep learning-based approaches, which offer strong performance in modeling spectral-spatial dependencies. DeepGUN [15] combines deep latent representation learning with vertex component analysis (VCA) for endmember extraction. Autoencoder-based models have also gained traction, with CNNAEU [16] being the first to introduce a convolutional autoencoder for hyperspectral unmixing. The LDVAE model [5] further advances this line by incorporating a latent Dirichlet distribution within a variational autoencoder (VAE) framework. The Dirichlet prior naturally satisfies the non-negativity and sum-to-one constraints of the unmixing problem while enabling a probabilistic interpretation of abundances. Building on this, SpACNN-LDVAE [6] replaces the MLP encoder with a CNN encoder, demonstrating improved performance by preserving spatial structure during training. Ghosh *et al.* introduced the first ViT-based model for hyperspectral unmixing, DeepTrans-HsU [3], which showed that vision transformers can outperform previous deep learning approaches with minimal architectural changes. UnDAT [2] employs a transformer encoder-decoder backbone along with spectral and spatial clustering modules to enhance unmixing performance, albeit at the cost of higher computational complexity. SSF-Net [17], an autoencoder-based model, introduces a spatial-spectral fusion module designed to capture local spatial heterogeneity and spectral diversity among endmembers more effectively.

3 Method

In hyperspectral imaging, we often represent each pixel $\mathbf{x} \in \mathbb{R}^C$ using

$$\mathbf{x} = \mathbf{E}\mathbf{z} + \mathbf{n}$$

where C refers to the number of channels, $\mathbf{E} \in \mathbb{R}^{C \times K}$ is the endmember matrix that collects K endmembers column-wise, $\mathbf{z} \in \mathbb{R}^{K \times 1}$ is the abundance vector representing the proportion of each endmember in pixel \mathbf{x} , and \mathbf{n} is additive

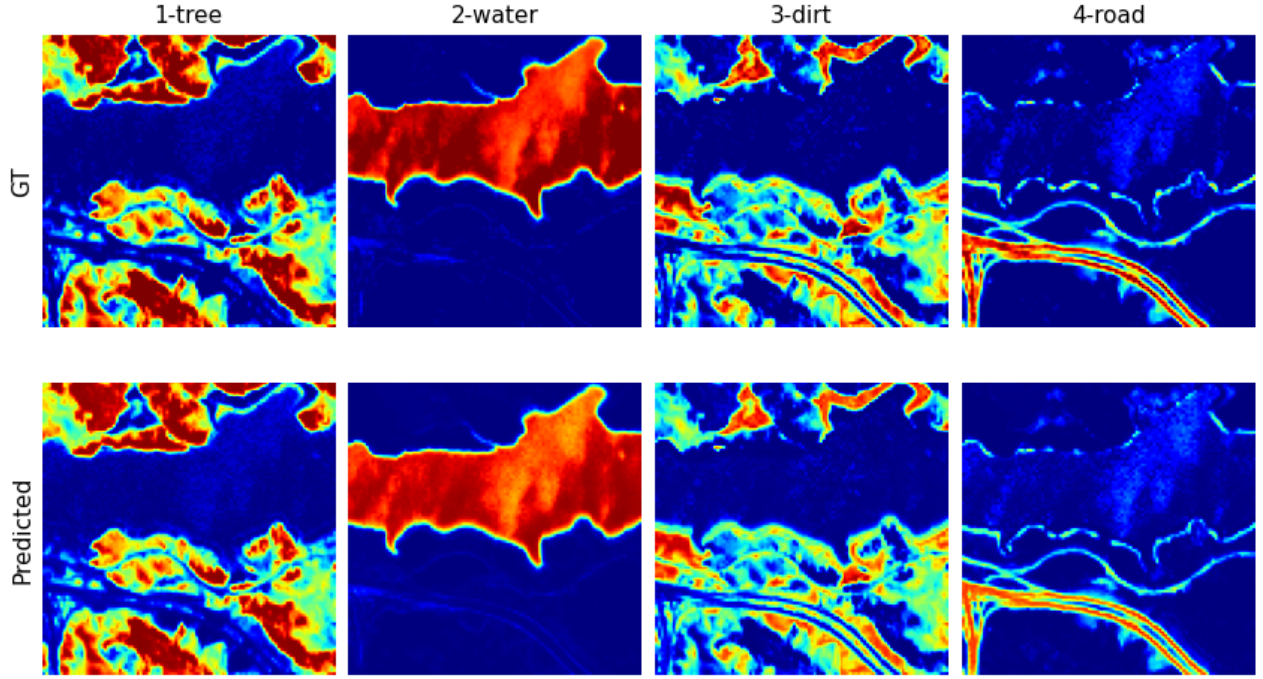


Figure 3: Ground Truth vs. Predicted endmember heat-maps for Jasper Ridge dataset. These plots visualize per-pixel abundances.

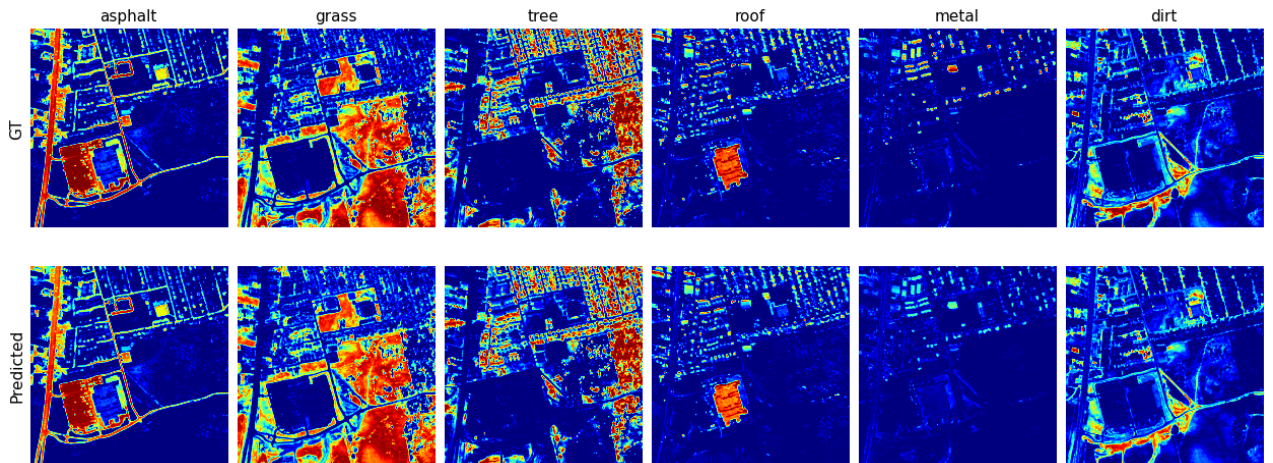


Figure 4: Ground Truth vs Predicted endmember heat-maps for HYDICE Urban dataset. These plots visualize per-pixel abundances.

noise. Here the abundance vector follows non-negative and sum-to-one constraints, i.e., $\forall_{k \in K} a_k \geq 0$ and $\sum_{k \in K} z_k = 1$. Unmixing computes abundances \mathbf{a} and endmembers matrix \mathbf{E} given a hyperspectral pixel \mathbf{x} . A straightforward extension is to consider a local patch $\mathcal{N}(\mathbf{x})$ around \mathbf{x} when performing unmixing.

3.1 LDVAE-T Model

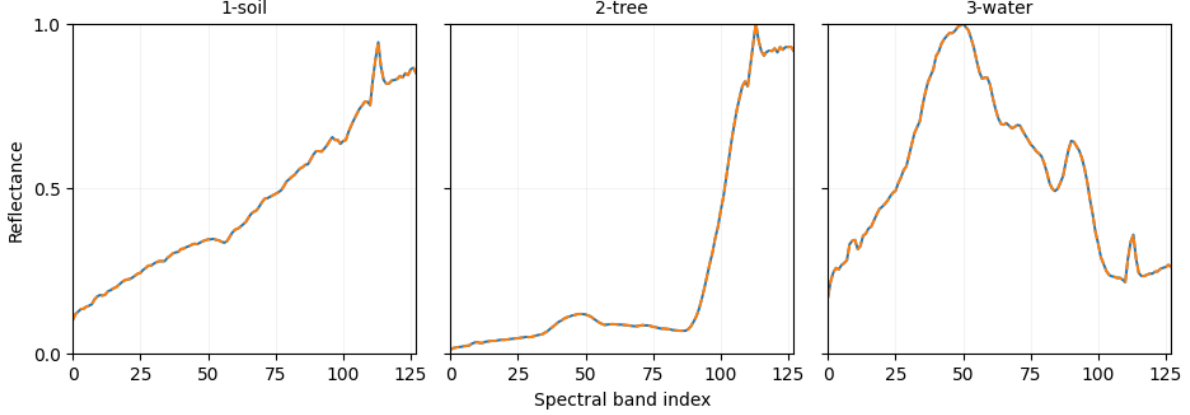


Figure 5: Extracted (orange) vs ground truth (blue) spectral signatures for endmembers in Samson dataset.

Table 1: SAD Scores for Endmember Extraction on Samson Dataset

	NMF-QMV [14]	CNNAEU [16]	LDVAE [5]	SpACNN-LDVAE [6]	DeepTrans-HsU [3]	UnDAT [2]	SSF-Net [17]	LDVAE-T
Endmember	2022	2021	2024	2024	2022	2023	2024	2025
soil	0.02326	0.07565	0.0959	0.2097	0.0128	0.01191	0.0092	0.000114 $\uparrow 98.76\%$
tree	0.06086	0.05440	1.2788	0.5347	0.0674	0.03775	0.0314	0.000202 $\uparrow 99.36\%$
water	1.45643	0.03642	0.4022	0.8233	0.0729	0.00813	0.0373	0.000213 $\uparrow 97.38\%$
average	0.51352	0.05549	0.5923	0.5525	0.0510	0.01926	0.0260	0.000177 $\uparrow 99.08\%$

Table 2: RMSE Scores for Abundance Estimation on Samson Dataset

	NMF-QMV [14]	CNNAEU [16]	LDVAE [5]	SpACNN-LDVAE [6]	DeepTrans-HsU [3]	UnDAT [2]	SSF-Net [17]	LDVAE-T
Endmember	2022	2021	2024	2024	2022	2023	2024	2025
soil	0.23298	0.3157	0.2609	0.2097	0.0712	0.04306	0.0511	0.03032 $\uparrow 29.59\%$
tree	0.24432	0.2911	0.3431	0.5347	0.0683	0.02854	0.0502	0.02401 $\uparrow 15.87\%$
water	0.37621	0.1552	0.3165	0.2098	0.0930	0.03128	0.0272	0.02249 $\uparrow 17.32\%$
average	0.28450	0.2540	0.3078	0.2412	0.0783	0.03429	0.0428	0.02561 $\uparrow 25.31\%$

The proposed *LDVAE-T* architecture consists of a transformer-based encoder that process a local hyperspectral image patch $\mathcal{P}(\mathbf{x})$ and predicts Dirichlet concentration parameters α (abundance prior) for pixel \mathbf{x} . Additionally, it emits a latent code for \mathbf{x} that the decoder uses to reconstruct the pixel’s endmember spectra. Recall that endmember spectra are shared by all pixels. Under the assumption that each endmember follows a Gaussian distribution, the decoder maps the latent code to the Gaussian parameters—mean μ_k and covariance Σ_k —for each endmember k . Next, we sample endmembers and mix them explicitly. For each endmember k , draw $\mathbf{e}_k \sim \mathcal{N}(\mu_k, \Sigma_k)$ and draw abundances $\mathbf{z} \sim \text{Dir}(\alpha)$. The pixel is then reconstructed by the mixture $\hat{\mathbf{x}} = \phi(\sum_k z_k \mathbf{e}_k)$.

Figure 1 illustrates the model. The tokenizer converts the local patch $\mathcal{P}(\mathbf{x})$ into S tokens $\mathbf{p}_s \in \mathbb{R}^d$ for $s \in \{1, \dots, S\}$. Learned positional encodings are added to each token,

$$\tilde{\mathbf{p}}_s = \mathbf{p}_s + \text{pos}_s.$$

The sequence $\{\tilde{\mathbf{p}}_s\}_{s=1}^S$ is passed to a transformer encoder with $L = 4$ layers; each layer uses multi-head self-attention with $H = 16$ heads followed by a position-wise feed-forward network (with residual connections and layer normalization). Self-attention enables the model to capture long-range dependencies across both spectral bands and

spatial locations within the patch, yielding a comprehensive representation of material mixtures. The encoder outputs $\{\mathbf{h}_s\}_{s=1}^S$ are max-pooled to construct the latent code for \mathbf{x}

$$\mathbf{x}_{\text{latent}} = \max_{s \in \{1, \dots, S\}} \mathbf{h}_s \text{ (element-wise max).}$$

The latent vector is passed through a softplus to produce the Dirichlet concentration parameters

$$\boldsymbol{\alpha} = \text{softplus}(\mathbf{W}\mathbf{x}_{\text{latent}} + \mathbf{b}) + \epsilon \mathbf{1}, \quad \boldsymbol{\alpha} \in \mathbb{R}_{>0}^K.$$

Here \mathbf{W} and \mathbf{b} are affine parameters for the soft-plus layer and small $\epsilon > 0$ ensures strictly positive entries. Sampling from $\text{Dir}(\boldsymbol{\alpha})$ yields an abundance vector \mathbf{z} that satisfies non-negativity and sum-to-one constraints.

The decoder comprises of two MLPs. First, MLP takes the latent vector $\mathbf{x}_{\text{latent}}$ and predicts endmembers means $\boldsymbol{\mu}_k$ and covariances Σ_k . The second MLP takes abundance vector \mathbf{z} and sampled endmembers $\mathbf{e}_k \sim \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k)$ and reconstruct pixel $\hat{\mathbf{x}}$. Similar to other approaches, we assume that the number of endmembers K is known *a priori*.

We employ the following losses during training. **Abundance loss** penalizes the divergence between predicted and ground truth abundances

$$\mathcal{L}_{\text{abundance}} = \text{MSE}(\hat{\mathbf{z}}, \mathbf{z})$$

where $\hat{\mathbf{z}}$ are the predicted abundances for pixel \mathbf{x} and \mathbf{z} are the ground-truth abundances.

Endmember bundles loss minimizes the KL divergence between the predicted endmember Gaussian and the ground-truth endmember bundle (estimated from *pure pixels*)

$$\mathcal{L}_{\text{endmembers}} = \sum_{k=1}^K \alpha_k \text{KL} \left(\mathcal{N}(\hat{\boldsymbol{\mu}}_k, \hat{\Sigma}_k) \parallel \mathcal{N}(\boldsymbol{\mu}_k, \Sigma_k) \right).$$

Where α_k is the predicted pre-sampled abundance estimation of endmember k . This mixture-weighted regularization constrains the loss calculation to endmembers present in the given pixel.

Lastly, the latent Dirichlet variational autoencoder optimizes

$$\mathcal{L}_{\text{ELBO}} = \mathbf{E}_{q_{\theta}(\mathbf{z}|\mathbf{x})} [\log p_{\phi}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\theta}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})).$$

with reconstruction term

$$\mathcal{L}_{\text{recon}} = \text{MSE}(\hat{\mathbf{x}}, \mathbf{x})$$

and Dirichlet prior $p(\mathbf{z}) = \text{Dir}(\boldsymbol{\alpha}^{\text{prior}})$. When $q_{\theta}(\mathbf{z}|\mathbf{x}) = \text{Dir}(\hat{\boldsymbol{\alpha}})$

$$\begin{aligned} \text{KL}(\text{Dir}(\hat{\boldsymbol{\alpha}}) \parallel \text{Dir}(\boldsymbol{\alpha}^{\text{prior}})) &= \sum \log \Gamma(\alpha_k^{\text{prior}}) \\ &\quad - \sum \log \Gamma(\hat{\alpha}_k) \\ &\quad + \sum (\hat{\alpha}_k - \alpha_k^{\text{prior}}) \frac{d}{dx} \ln \Gamma(\hat{\alpha}_k). \end{aligned}$$

Here $\hat{\boldsymbol{\alpha}}$ is the concentration parameter of the estimated Dirichlet distribution and $\boldsymbol{\alpha}^{\text{prior}}$ is the concentration parameter of the Dirichlet prior. $\Gamma(\cdot)$ is the Gamma function. Thus the overall loss is

$$\mathcal{L} = \mathcal{L}_{\text{ELBO}} + \lambda_{\text{abundances}} \mathcal{L}_{\text{abundances}} + \lambda_{\text{endmembers}} \mathcal{L}_{\text{endmembers}}$$

with $\lambda_{\text{abundances}}$ and $\lambda_{\text{endmembers}}$ control the relative strengths of abundances and endmember terms. We set $\lambda_{\text{abundances}}$ to 1, and apply a heavy anneal to $\lambda_{\text{endmembers}}$ for experiments reported here. The annealing begins at 1×10^{-6} and progresses towards 1.0 over 80000 epochs (which is never reached). This is done to keep $\mathcal{L}_{\text{endmembers}}$ from significantly outweighing the other terms by a magnitude of 1×10^6 at the start of training.

4 Experimental Setup

We evaluate our model on three widely used hyperspectral unmixing benchmarks: Samson [18], Jasper Ridge [19], and HYDICE Urban [20]. The Samson dataset contains a 95×95 hyperspectral image with 156 spectral bands and three ground-truth endmembers: *Soil*, *Tree*, and *Water*. Jasper Ridge comprises a 100×100 image with 198 spectral bands and four ground-truth endmembers: *Tree*, *Water*, *Dirt*, and *Road*. HYDICE Urban has 307×307 pixels with 162 spectral bands and is available in three variants containing four, five, or six endmembers. We use the six-endmember variant: *Asphalt Road*, *Grass*, *Tree*, *Roof*, *Metal*, and *Dirt*. All three datasets provide per-pixel ground-truth abundance maps.

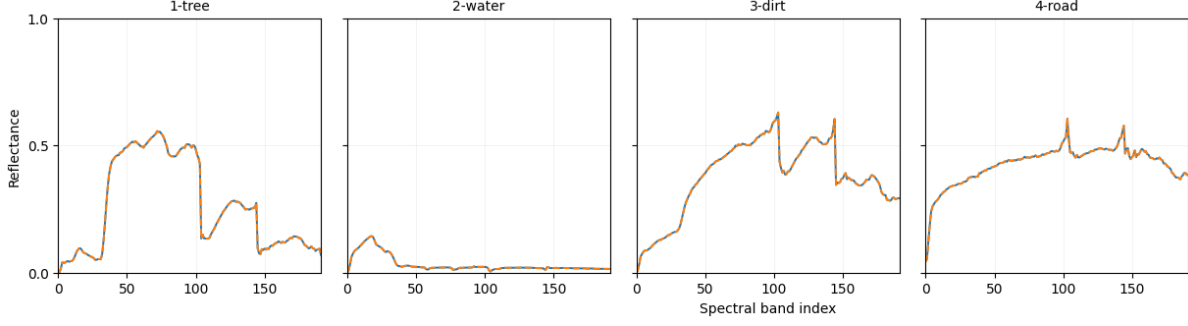


Figure 6: Extracted (orange) vs ground truth (blue) spectral signatures for endmembers in Jasper Ridge dataset.

Table 3: SAD Scores for Endmember Extraction on Jasper Ridge Dataset

Endmember	NMF-QMV [14] 2022	CNNAEU [16] 2021	LDVAE [5] 2024	SpACNN-LDVAE [6] 2024	DeepTrans-HsU [3] 2022	UnDAT [2] 2023	SSF-Net [17] 2024	LDVAE-T 2025
tree	0.05016	0.3104	-	-	-	0.04519	0.0781	0.000255 $\uparrow 99.44\%$
water	0.28387	0.6082	-	-	-	0.02811	0.0293	0.001269 $\uparrow 95.49\%$
soil	0.17326	0.3381	-	-	-	0.09074	0.0484	0.000270 $\uparrow 99.44\%$
road	1.46974	0.0519	-	-	-	0.03306	0.0238	0.000288 $\uparrow 98.79\%$
average	0.49426	0.3271	-	-	-	0.04927	0.0449	0.000520 $\uparrow 98.84\%$

Table 4: RMSE Scores for Abundance Estimation on Jasper Ridge Dataset

Endmember	NMF-QMV [14] 2022	CNNAEU [16] 2021	LDVAE [5] 2024	SpACNN-LDVAE [6] 2024	DeepTrans-HsU [3] 2022	UnDAT [2] 2023	SSF-Net [17] 2024	LDVAE-T 2025
tree	0.12528	0.3169	-	-	-	0.06031	0.0721	0.02609 $\uparrow 56.74\%$
water	0.20375	0.2118	-	-	-	0.04211	0.0761	0.02892 $\uparrow 31.32\%$
soil	0.14647	0.2978	-	-	-	0.07007	0.0930	0.03855 $\uparrow 44.98\%$
road	0.18446	0.2043	-	-	-	0.07792	0.0782	0.03819 $\uparrow 50.99\%$
average	0.16499	0.2577	-	-	-	0.06260	0.0798	0.03294 $\uparrow 47.38\%$

These datasets provide a single “true” spectrum for each endmember. In practice, endmember spectra are not unique: the same material can exhibit measurable spectral variability under different illumination, viewing geometry, sensor characteristics, and environmental conditions. We first identify high-purity pixels using the Pixel Purity Index (PPI)—pixels dominated by a single endmember—and then use their spectra to initialize the spectral distribution of each endmember. We model each endmember’s spectral variability with a (multivariate) Gaussian distribution. We refer to these as endmember bundles, that is, each endmember is represented by a distribution over spectra rather than a single spectral signature.

We adopt a 20/80 train–test split for training and evaluation. The model is trained in a supervised setting for 1000 epochs using a batch size of 128 and the Adam optimizer with a learning rate of 2×10^{-4} . During training, zero-padding is applied to the image boundaries to maintain spatial consistency in patch extraction.

4.1 Evaluation Metrics

We capture endmember reconstruction using Spectral Angle Distance (SAD), which quantifies the discrepancy between predicted and ground-truth endmember spectra by measuring the angle between their spectral vectors in a high-dimensional space, making it invariant to per-pixel intensity (scaling) changes:

$$\text{SAD}(\hat{\mathbf{e}}, \mathbf{e}) = \cos^{-1} \left(\frac{\hat{\mathbf{e}}_k^\top \mathbf{e}_k}{\|\hat{\mathbf{e}}_k\| \|\mathbf{e}_k\|} \right),$$

where $\hat{\mathbf{e}}_k$ is the predicted endmember distribution’s mean and \mathbf{e}_k is the mean spectrum of the corresponding GT endmember bundle. Lower SAD values indicate closer alignment between the predicted and reference spectra. We use SAD to assess endmember-reconstruction accuracy for each material class in the dataset.

We evaluate abundance-estimation accuracy using the Root Mean Squared Error (RMSE), which measures the average deviation between predicted and ground-truth abundances across pixels:

$$\text{RMSE}(\hat{\mathbf{z}}, \mathbf{z}) = \sqrt{\frac{1}{N} \sum_{n=1}^N \|\hat{\mathbf{z}}_n - \mathbf{z}_n\|^2}$$

where $\hat{\mathbf{z}}$ and \mathbf{z} refers to the predicted and the ground truth abundances. N is the number of pixels. A lower RMSE value indicates better abundance estimation.

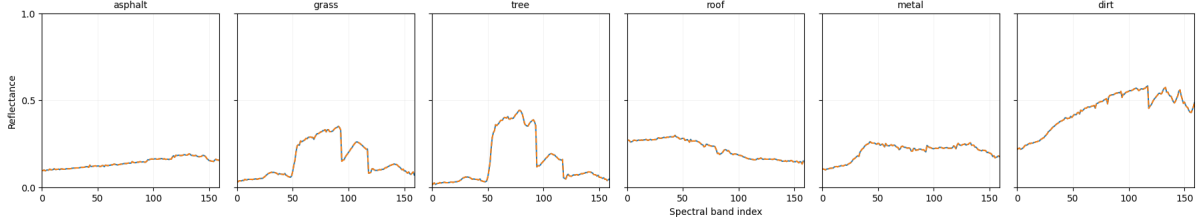


Figure 7: Extracted (orange) vs ground truth (blue) spectral signatures for endmembers in HYDICE Urban dataset.

Table 5: SAD Scores for Endmember Extraction on HYDICE Urban Dataset. *SSF-Net used a five endmember version of HYDICE Urban for their experiments.

Endmember	SGSNMF [11] 2017	SSWNMF [10] 2022	LDVAE [5] 2024	SpACNN-LDVAE [6] 2024	SSF-Net* [17] 2024	LDVAE-T 2025
Asphalt road	0.0841	0.0782	0.4262	0.2786	0.0629	0.000216 ↑99.66%
Grass	0.1516	0.1490	0.3323	0.1936	0.0411	0.000145 ↑99.65%
Tree	0.1199	0.1173	0.3177	0.4411	0.0850	0.000155 ↑99.82%
Roof	0.0731	0.0713	0.4393	0.4502	0.0487	0.000208 ↑99.57%
Metal	0.1250	0.1241	0.7004	0.3241	-	0.000338 ↑99.73%
Dirt	0.0859	0.0802	0.2806	0.2026	0.1065	0.000088 ↑99.89%
Average	0.1060	0.1034	0.4161	0.3151	0.0688	0.000192 ↑99.72%

Table 6: RMSE Scores for Abundance Estimation on HYDICE Urban Dataset. *SSF-Net used a five endmember version of HYDICE Urban for their experiments.

Endmember	SGSNMF [11] 2017	SSWNMF [10] 2022	LDVAE [5] 2024	SpACNN-LDVAE [6] 2024	SSF-Net* [17] 2024	LDVAE-T 2025
Asphalt road	-	-	0.2889	0.1566	0.1578	0.03538 ↑77.41%
Grass	-	-	0.1832	0.1977	0.1416	0.03083 ↑78.23%
Tree	-	-	0.1737	0.1632	0.1179	0.02576 ↑78.15%
Roof	-	-	0.1250	0.1283	0.0909	0.02476 ↑72.76%
Metal	-	-	0.2599	0.0992	-	0.03830 ↑61.39%
Dirt	-	-	0.1334	0.1894	0.1192	0.03128 ↑73.76%
Average	-	-	0.1840	0.1558	0.1256	0.03105 ↑75.28%

5 Results

We evaluate the proposed model on two key hyperspectral unmixing tasks: (1) abundance estimation and (2) endmember extraction.

Abundance Estimation. Tables 2, 4, and 6 report Root Mean Squared Error (RMSE) for abundance estimation on the Samson, Jasper Ridge, and HYDICE Urban datasets. Our model consistently achieves lower RMSE across all endmembers and datasets, outperforming state-of-the-art baselines. This performance is attributable to the Dirichlet prior imposed in the latent space, which naturally enforces the sum-to-one and non-negativity constraints expected of abundance vectors, yielding stable and physically meaningful estimates. By contrast, methods based on unconstrained regression can violate these constraints or require *ad hoc* normalization. Figures 2, 3, and 4 show heatmaps of the predicted abundances vs. ground truth abundances for the Samson, Jasper Ridge, and HYDICE Urban datasets.

Endmember Extraction. Tables 1, 3, and 5 present Spectral Angle Distance (SAD) for the Samson, Jasper Ridge, and HYDICE Urban datasets. Our model also outperforms state-of-the-art methods by a significant margin across all datasets and endmembers, achieving uniformly lower SAD. This improvement can be attributed to the model’s use of endmembers bundles to more accurately estimate the distribution of endmembers that exist within a real world scene. Figures 5, 6, and 7 show extracted vs. ground truth spectra for Samson, Jasper Ridge, and HYDICE Urban datasets.

6 Conclusion and Future Work

We introduce the *Latent Dirichlet Transformer Variational Autoencoder* (LDVAE-T), a framework for hyperspectral unmixing that couples transformer-based modeling with a Dirichlet latent structure. By imposing a Dirichlet prior in the latent space, the model naturally enforces the sum-to-one and non-negativity constraints required for physically plausible abundance estimates. Beyond this probabilistic scaffold, LDVAE-T contributes a decoder based on *bundled endmembers*: instead of treating each material as a single, fixed prototype, the decoder predicts for each patch a distributional prototype comprising a mean spectrum and a structured covariance defined over spectral segments. Reconstructions are formed by mixing these bundles with Dirichlet-distributed abundances, enabling the network to capture intrinsic intra-class variability while preserving interpretability. Experiments on three standard benchmarks—*Samson*, *Jasper Ridge*, and *HYDICE Urban*—show that LDVAE-T delivers state-of-the-art performance in both endmember extraction and abundance estimation, as measured by spectral angle distance and root mean squared error, respectively.

The empirical gains of LDVAE-T stem from the combining of three ingredients: 1) a *physically meaningful latent space* via the Dirichlet prior, which allows for a VAE framework with a distribution that works within the constraints of a mixture; 2) *distributional endmembers* that explicitly model spectral variability present in real world scenes; and 3) *transformer-based encoding* that leverages long-range spectral–spatial dependencies often missed by MLP or CNN encoders. We find that Segment-wise covariance parameterization strikes a balance between flexibility (capturing correlated variation within bands) and tractability (avoiding full dense covariances over all wavelengths). Collectively, these elements improve both abundance estimation and the stability of extracted endmember signatures.

References

- [1] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale,” 2021. Copyright - © 2021. This work is published under <http://arxiv.org/licenses/nonexclusive-distrib/1.0/> (the “License”). Notwithstanding the ProQuest Terms and Conditions, you may use this content in accordance with the terms of the License; Last updated - 2024-10-16.
- [2] Y. Duan, X. Xu, T. Li, B. Pan, and Z. Shi, “Undat: Double-aware transformer for hyperspectral unmixing,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 61, pp. 1–12, 2023.
- [3] P. Ghosh, S. K. Roy, B. Koirla, B. Rasti, and P. Scheunders, “Hyperspectral unmixing using transformer network,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–16, 2022.
- [4] Z. Yang, M. Xu, S. Liu, H. Sheng, and J. Wan, “Ust-net: A u-shaped transformer network using shifted windows for hyperspectral unmixing,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 61, pp. 1–15, 2023.
- [5] K. Mantripragada and F. Z. Qureshi, “Hyperspectral pixel unmixing with latent dirichlet variational autoencoder,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 62, pp. 1–12, 2024.
- [6] S. Chitnis, K. Mantripragada, and F. Z. Qureshi, “Spacnn-ldvae: Spatial attention convolutional latent dirichlet variational autoencoder for hyperspectral pixel unmixing,” in *IGARSS 2024 - 2024 IEEE International Geoscience and Remote Sensing Symposium*, pp. 7714–7719, 2024.
- [7] L. Sun and P. G. Lucey, “Unmixing mineral abundance and mg# with radiative transfer theory: Modeling and applications,” *Journal of Geophysical Research: Planets*, vol. 126, no. 2, p. e2020JE006691, 2021. e2020JE006691 2020JE006691.

- [8] J. Janiczek, P. Thaker, G. Dasarathy, C. S. Edwards, P. Christensen, and S. Jayasuriya, “Differentiable programming for hyperspectral unmixing using a physics-based dispersion model,” 2020. Copyright - © 2020. This work is published under <http://arxiv.org/licenses/nonexclusive-distrib/1.0/> (the “License”). Notwithstanding the ProQuest Terms and Conditions, you may use this content in accordance with the terms of the License; Last updated - 2020-07-15.
- [9] L. Drumetz, J. Chanussot, and C. Jutten, “Spectral unmixing: A derivation of the extended linear mixing model from the hapke model,” *IEEE Geoscience and Remote Sensing Letters*, vol. 17, no. 11, pp. 1866–1870, 2020.
- [10] S. Zhang, G. Zhang, F. Li, C. Deng, S. Wang, A. Plaza, and J. Li, “Spectral-spatial hyperspectral unmixing using nonnegative matrix factorization,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–13, 2022.
- [11] X. Wang, Y. Zhong, L. Zhang, and Y. Xu, “Spatial group sparsity regularized nonnegative matrix factorization for hyperspectral unmixing,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 11, pp. 6287–6304, 2017.
- [12] X. Lu, H. Wu, Y. Yuan, P. Yan, and X. Li, “Manifold regularized sparse nmf for hyperspectral unmixing,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 51, no. 5, pp. 2815–2826, 2013.
- [13] W. He, H. Zhang, and L. Zhang, “Total variation regularized reweighted sparse nonnegative matrix factorization for hyperspectral unmixing,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 7, pp. 3909–3921, 2017.
- [14] M. Zhao, T. Gao, J. Chen, and W. Chen, “Hyperspectral unmixing via nonnegative matrix factorization with handcrafted and learned priors,” *IEEE Geoscience and Remote Sensing Letters*, vol. 19, pp. 1–5, 2022.
- [15] R. A. Borsoi, T. Imbiriba, and J. C. M. Bermudez, “Deep generative endmember modeling: An application to unsupervised spectral unmixing,” *IEEE Transactions on Computational Imaging*, vol. 6, pp. 374–384, 2020.
- [16] B. Palsson, M. O. Ulfarsson, and J. R. Sveinsson, “Convolutional autoencoder for spectral–spatial hyperspectral unmixing,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 1, pp. 535–549, 2021.
- [17] B. Wang, H. Yao, D. Song, J. Zhang, and H. Gao, “Ssf-net: A spatial–spectral features integrated autoencoder network for hyperspectral unmixing,” *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 17, pp. 1781–1794, 2024.
- [18] H. I. Laboratory, “Samson hyperspectral dataset,” 2011. Publicly available dataset commonly used for hyperspectral unmixing studies.
- [19] N. J. P. Laboratory, “AVIRIS jasper ridge hyperspectral dataset.” <https://aviris.jpl.nasa.gov/>, 1995. Accessed: 2025-04-30.
- [20] U. A. C. of Engineers, “Hydice urban hyperspectral dataset,” 1995. Captured by the Hyperspectral Digital Imagery Collection Experiment (HYDICE) sensor. Available via academic requests.