

# Compression of Solar Spectroscopic Observations: Case Study of Mg II k Spectral Line Profiles Observed by NASA's IRIS Satellite

Viacheslav Sadykov<sup>1</sup>, Irina Kitiashvili<sup>2</sup>, Alberto Sainz Dalda<sup>3</sup>, Vincent Oria<sup>4</sup>, Alexander Kosovichev<sup>4</sup>, Egor Illarionov<sup>5</sup>

*<sup>1</sup>Georgia State University*

*<sup>2</sup>NASA Ames Research Center*

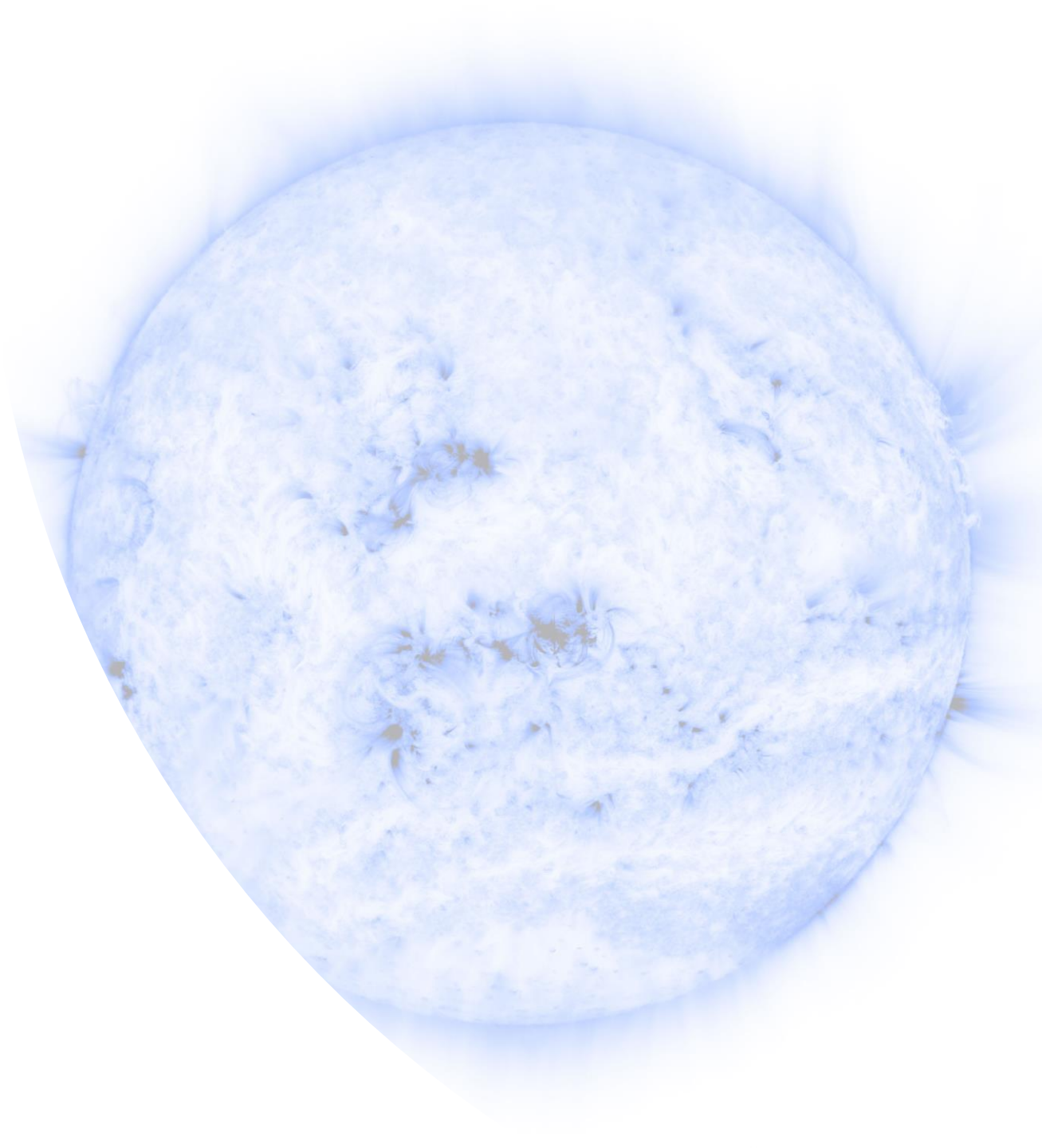
*<sup>3</sup>Lockheed Martin Solar and Astrophysics Laboratory*

*<sup>4</sup>New Jersey Institute of Technology*

*<sup>5</sup>Moscow State University*

# Outline

- Introduction to autoencoders
- Description of the IRIS spectral data and training procedure
- Results of the line profile compression
- Interpretation of the derived deep features
- Conclusions

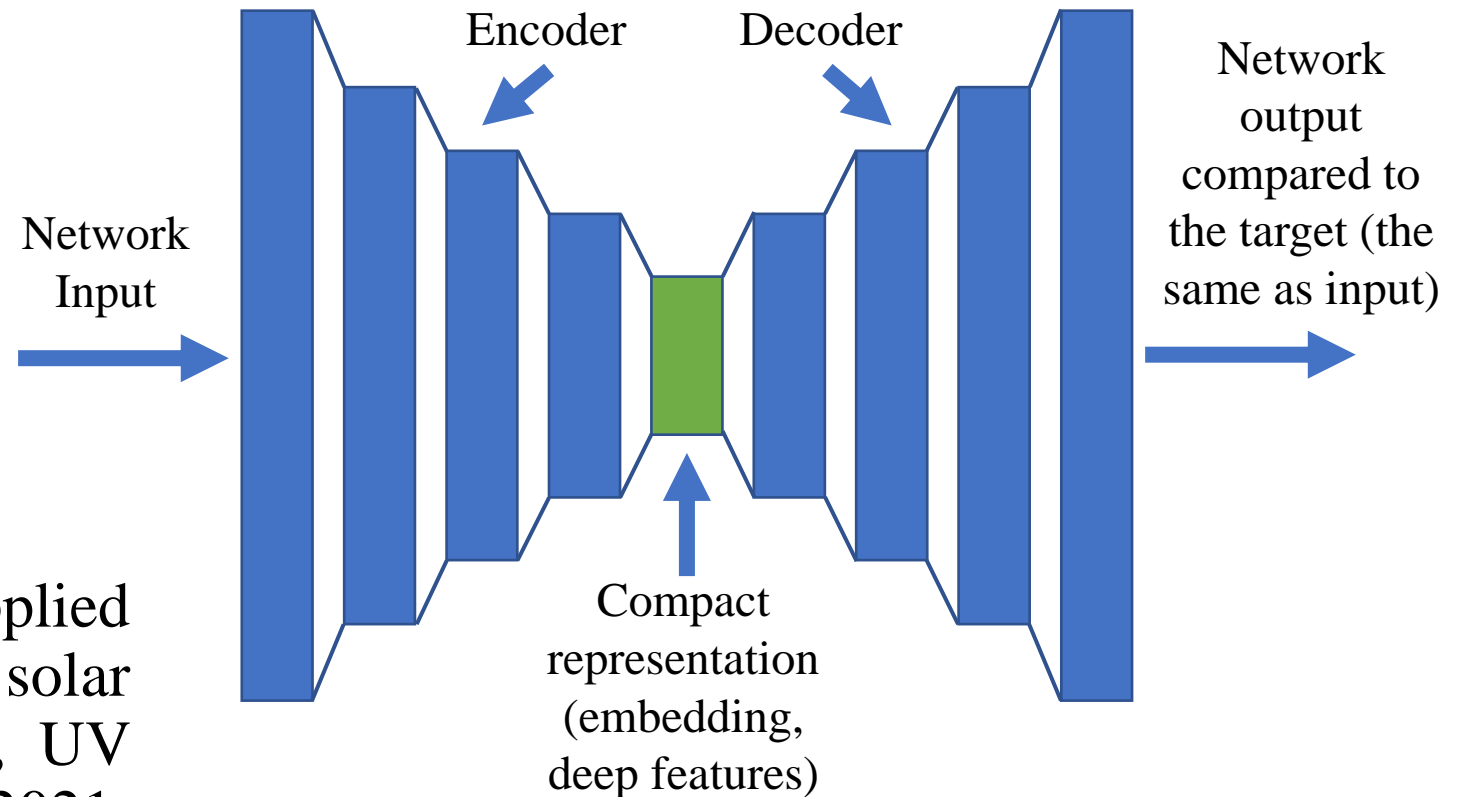


# Extraction of a compact representation with autoencoders

An autoencoder is a type of artificial neural network used to learn an efficient representation of the data (features) in an unsupervised manner. Typical applications of the autoencoders include:

- Noise reduction
- Extraction of deep features
- Dimensionality reduction
- Data compression
- Anomaly detection

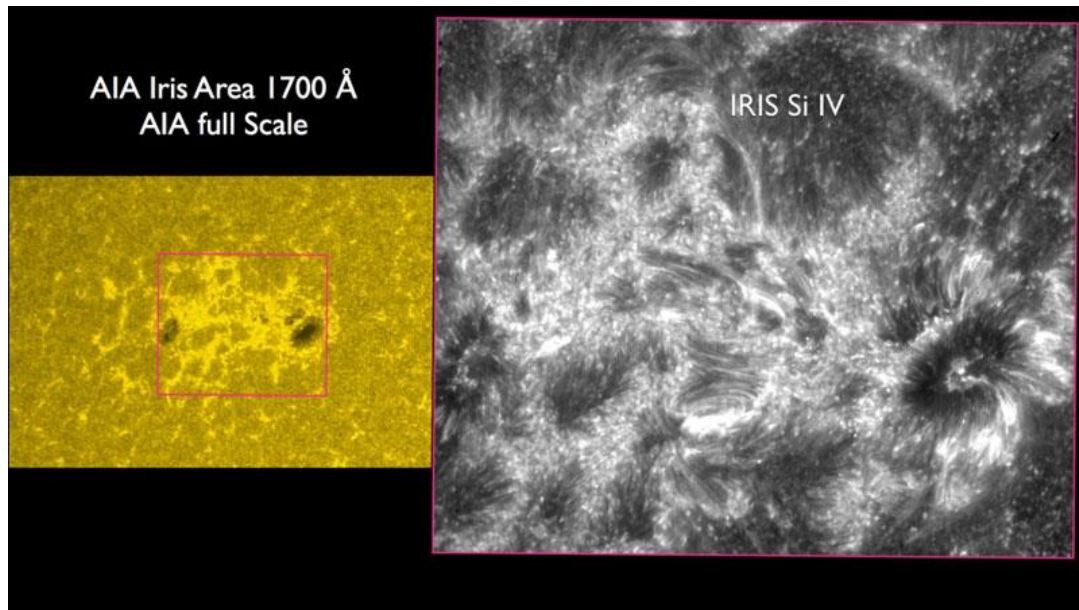
Autoencoders have already been applied for extraction of deep features from solar magnetograms (Chen et al. 2019), UV spectral lines (Panos et al. 2021, Sadykov et al. 2021), and more.



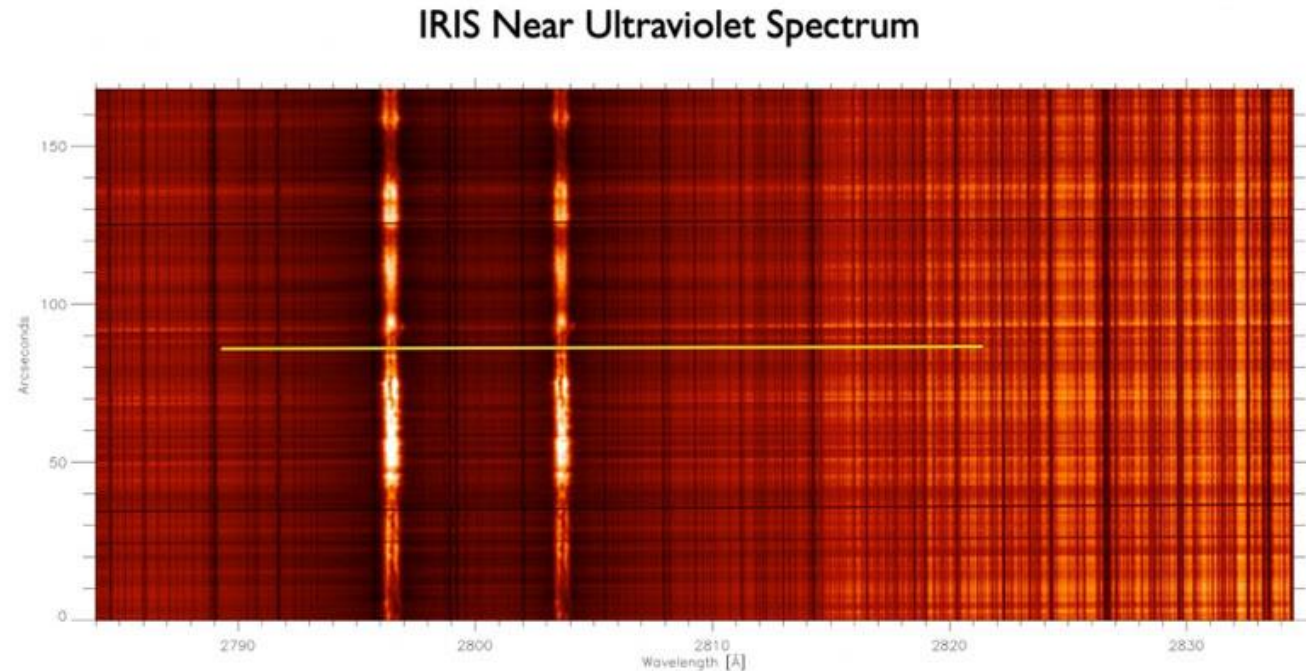
*Schematic illustration of an autoencoder*

# Interface Region Imaging Spectrograph (IRIS)

- IRIS is a NASA small explorer mission carrying an ultraviolet telescope combined with an imaging spectrograph, launched in 2013
- IRIS obtains simultaneous high-resolution observations of slit-jaw images (0.33'' angular resolution) and UV spectra (26-53 mÅ spectral resolution depending on the line; a variety of lines cover the photosphere, chromosphere, transition region, and corona)
- In this study, we focused on the Mg II k 2796 Å spectral line



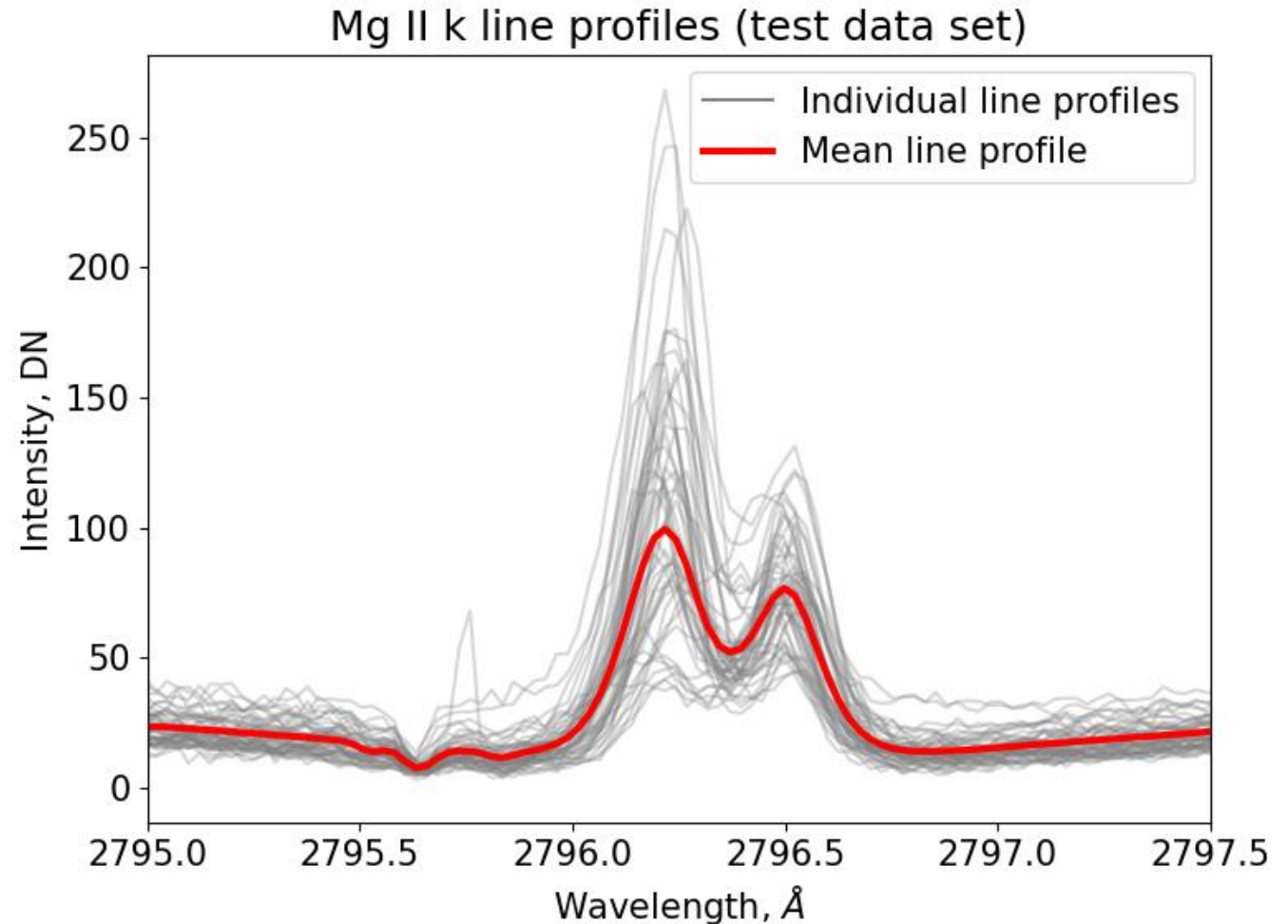
*Comparing resolutions of SDO/AIA and IRIS. Credit: NASA*



*An example of IRIS Mg II spectral data. Credit: NASA*

# Data set cleaning and train-test separation

- We utilized IRIS observations of the quiet Sun taken on April 20, 2020, from 08:32:00 UT - 09:56:00 UT at the center of the solar disk.
- The observations were made in the sit-and-stare mode; more than 300,000 individual Mg II k line profiles were obtained.
- After cleaning and normalizing the data set, we separated it into train and test subsets based on time.



*Illustration of Mg II k line profiles from the test data set*

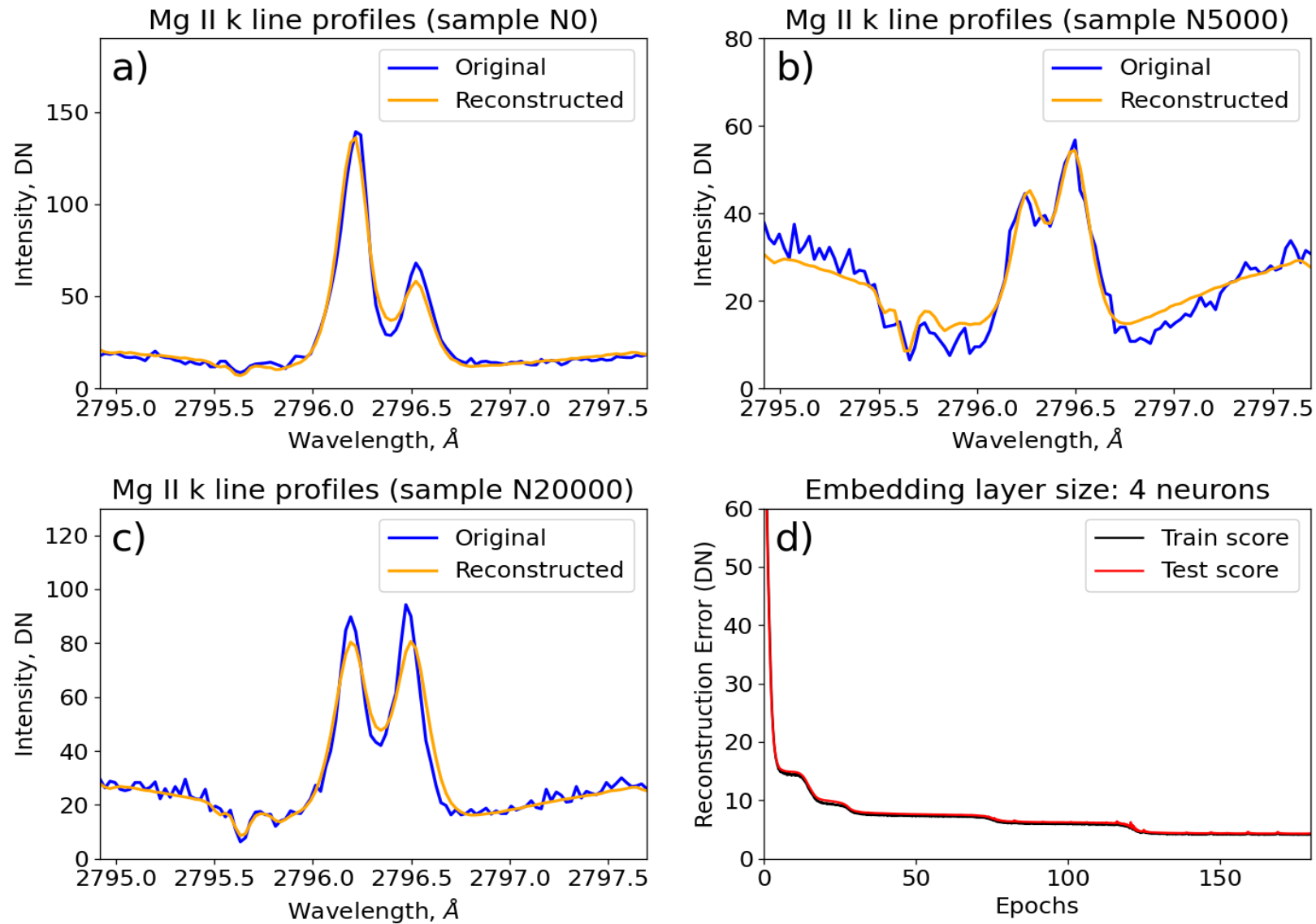
# Autoencoder architecture and training strategy

We utilized a fully-connected autoencoder for line profile compression:

- The autoencoder architecture (number of neurons in subsequent layers) was 110 - 64 - 32 - 16 -  $n_{\text{emb}}$  - 16 - 32 - 64 - 110.
- The number of neurons in the embedding layer,  $n_{\text{emb}}$ , varied from 1 to 15. The optimal number was an object of study.
- The neurons had Rectified Linear Unit (ReLU) activation functions.

The training strategy was the following:

- The loss function was the Mean Square Error (MSE, described later).
- An early stopping criterion (improvement of less than 0.5% on the test data for two subsequent epochs) was applied to prevent the network from overfitting.
- Several optimizers (adam and SGD) and learning rates were tested.
- For each  $n_{\text{emb}}$  the training process was repeated five times to estimate how the weight initialization affects the results



*Examples of the original (blue) and reconstructed (orange) Mg II k line profiles from the test subset for the case of an embedding layer size of 4 neurons. Lower right panel illustrates the corresponding autoencoder training process (mean squared error, MSE, as a function of the epoch number).*

# Evaluation of the line profile reconstruction

The performance of the autoencoder was evaluated based on the Mean Square Error (MSE) measure calculated as:

$$MSE = \sqrt{\frac{1}{k} \sum_{i=1}^k (I_i - I'_i)^2}$$

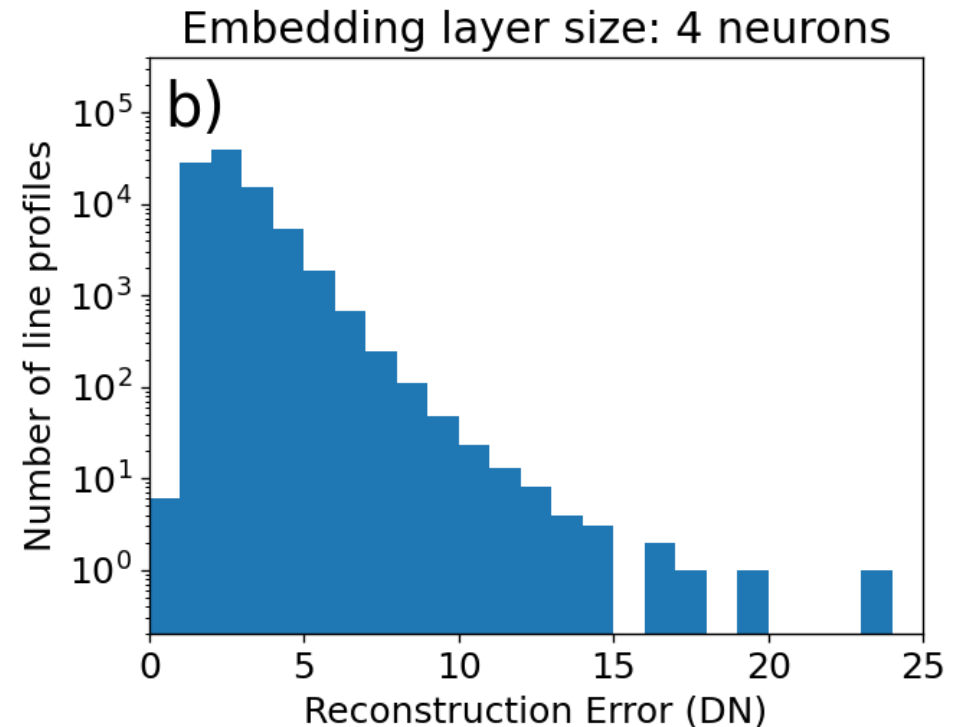
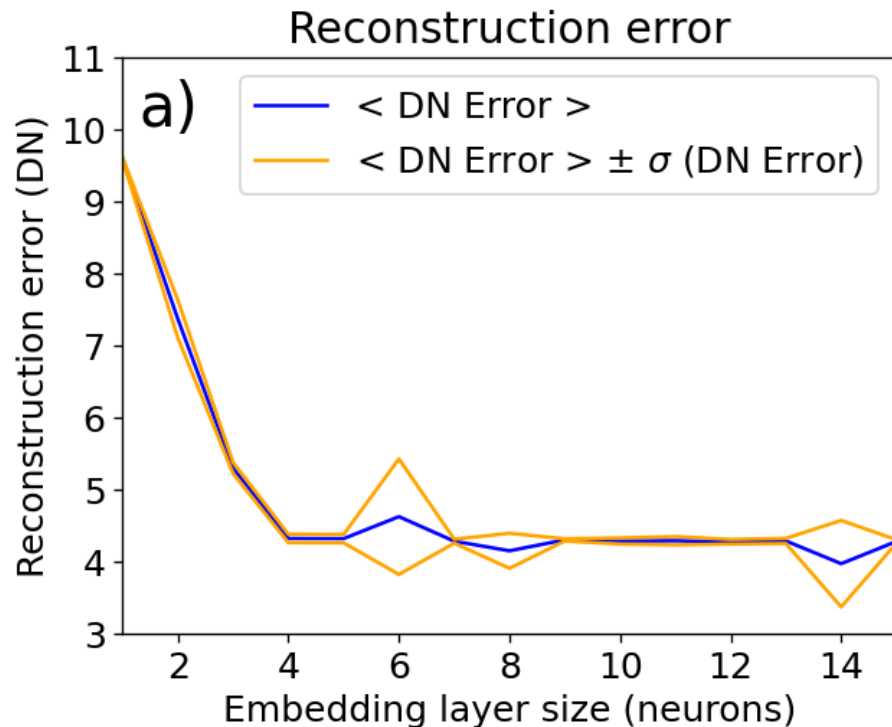
Here  $k$  is the number of wavelength points, and  $I_i$  and  $I'_i$  correspond to the original and reconstructed intensities of the line profile. In addition, we compared the first 10 statistical moments of the original and reconstructed line profiles computed as:

$$\mu_n = \left( \frac{\int_{-\infty}^{\infty} (\lambda - \lambda_0)^n I(\lambda) d\lambda}{\int_{-\infty}^{\infty} I(\lambda) d\lambda} \right)^{1/n},$$
$$\lambda_0 = \frac{\int_{-\infty}^{\infty} \lambda I(\lambda) d\lambda}{\int_{-\infty}^{\infty} I(\lambda) d\lambda}$$

Here  $I(\lambda)$  indicates the original or reconstructed intensity at wavelength  $\lambda$ , and  $n$  is the statistical moment number.

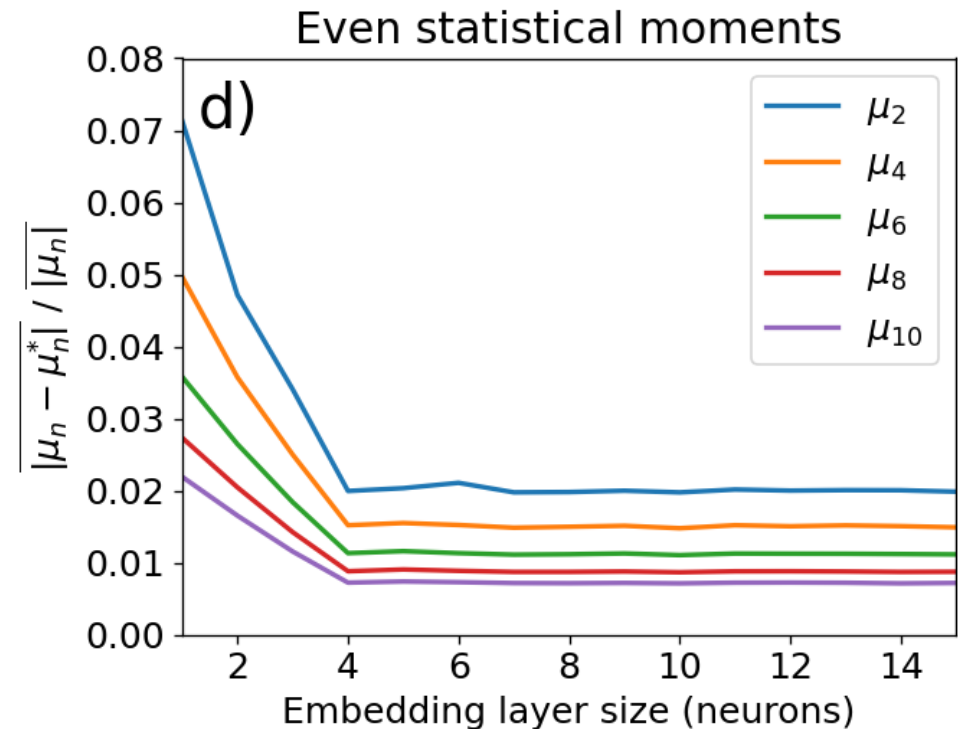
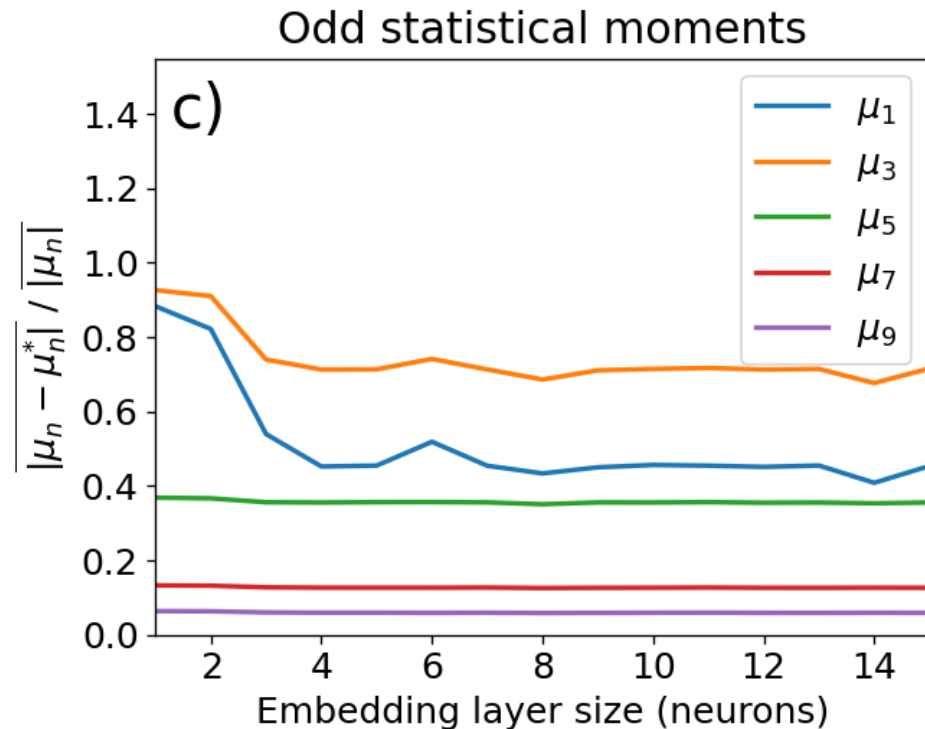
# Results: MSE reconstruction

- The MSE decreased sharply for  $n_{\text{emb}} \leq 4$  and barely dropped for  $n_{\text{emb}} > 4$ . One can state that there is almost no further progress in line compression for  $n_{\text{emb}} > 4$ .
- For most of the line profiles, the average deviation of the reconstructed intensity from the original values was around 3 DN, which is comparable with the intensity variations of the line continuum signal.



# Results: statistical moments

- The error of reconstruction of even moments improved for  $n_{\text{emb}} \leq 4$  and stopped decreasing for  $n_{\text{emb}} > 4$  (same as for MSE). Even moments have an average relative error of the reconstruction of  $\leq 2\%$ ; the error decreases for higher moments.
- Reconstruction of the odd moments had much larger relative errors and did not improve so sharply with increasing  $n_{\text{emb}}$ .



# Interpretation of Deep Features

To better understand the meaning of the parameters that the autoencoder learned for the  $n_{\text{emb}} = 4$  case, we conducted the following experiment on the embedding space:

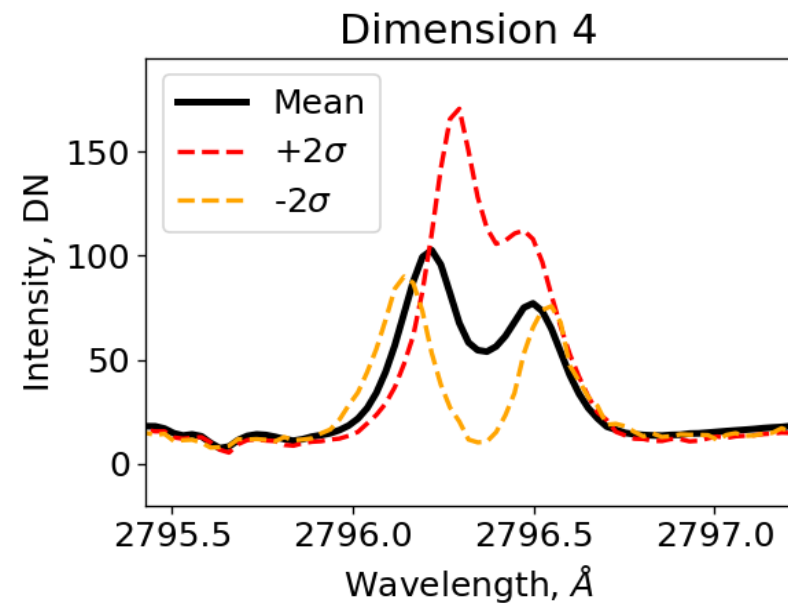
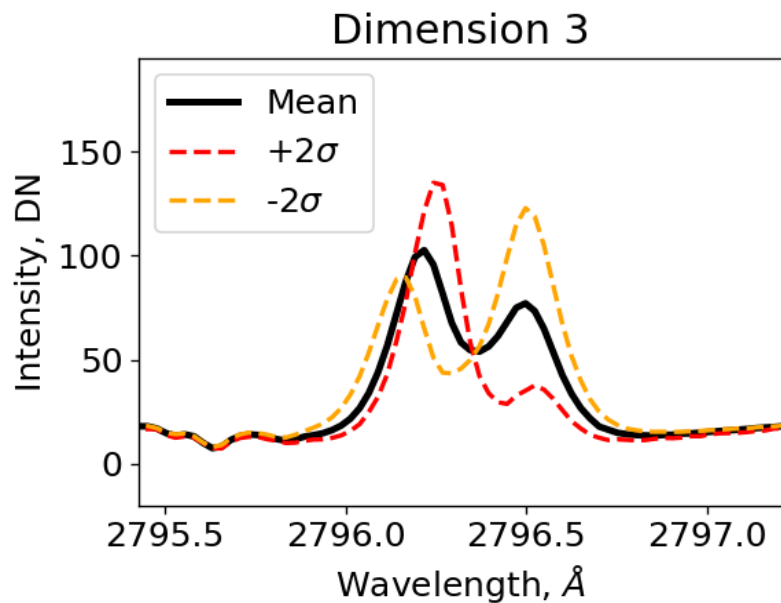
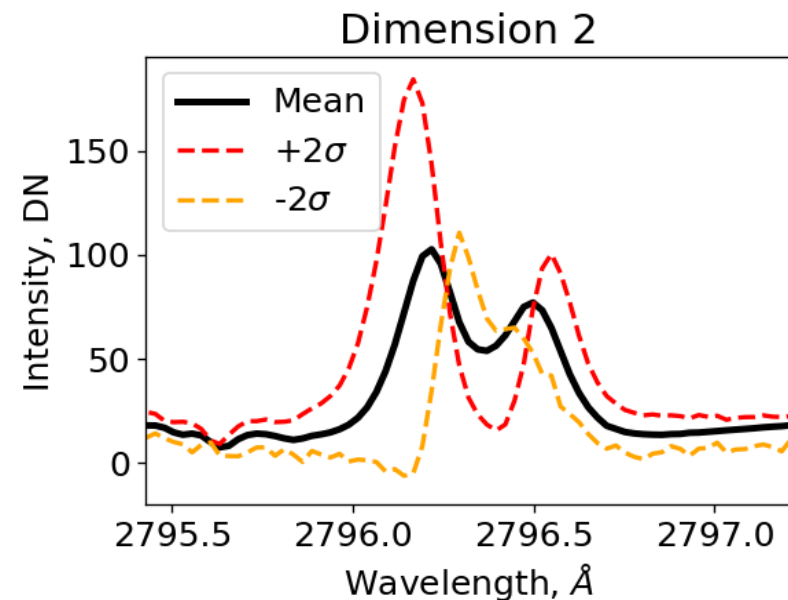
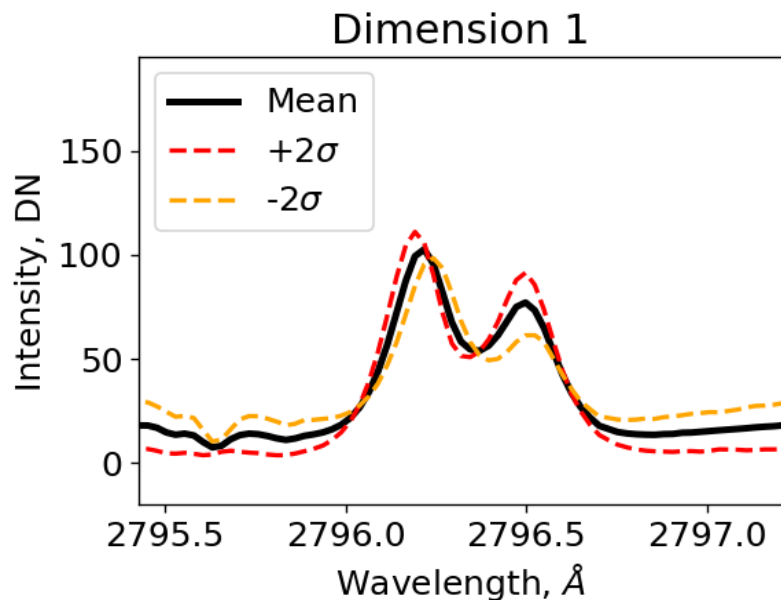
- We calculated the mean values and the standard deviations of the learned deep features across the test data set;
- We varied one of the four deep features in the range of the mean  $\pm$  two standard deviations and propagated these through the decoder.

These results demonstrate how variations in the embedding space affect the line profile shapes and could potentially help interpret the learned deep features.

# Interpretation of Deep Features

The meaning of the deep features learned by the autoencoder can be intuitively interpreted as follows:

- Dimension (feature) 3 is responsible for the line profile asymmetry
- Dimension 4 controls the central-reversal feature depth
- Dimension 2 affects the line profile width
- Dimension 1 seems to affect the line profile continuum



# Summary of the results

- It is possible to compress the data more than 27 times while having a reconstruction error somewhat comparable to the variations of the measurements in the line continuum;
- The average error of reconstruction of the MSE and even statistical moments of the line profiles decreases sharply for the  $n_{\text{emb}} \leq 4$  and barely drops for  $n_{\text{emb}} > 4$ . Reconstruction of the odd statistical line moments does not demonstrate such dependence on the dimensionality of the embedding layer;
- Two occasional improvements of the MSE were observed for  $n_{\text{emb}} > 4$  indicating complications in obtaining better embedding;
- The features learned for the  $n_{\text{emb}} = 4$  case can be supported with an intuitively-understandable interpretation when variations in the embedding space are considered.



Thank You for  
Your Attention!