

Technische Universität Berlin

Faculty of Electrical Engineering and Computer Science
Dept. of Computer Engineering and Microelectronics
Remote Sensing Image Analysis Group



Learning-Based Hyperspectral Image Compression Using A Spatio-Spectral Approach

Master of Science in Computer Science

August, 2023

Niklas Sprengel

Matriculation Number: 380009

Supervisor: Prof. Dr. Begüm Demir

Advisor: Martin Fuchs

Sworn Affidavit

I hereby declare that the thesis submitted is my own, unaided work, completed without any unpermitted external help. Only the sources and resources listed were used.

Berlin, Date

.....

Niklas Sprengel

Abstract

Learning-based hyperspectral image compression is an increasingly important field of study in the field of remote sensing. This thesis proposes the Combined Model, a novel, highly flexible architecture to perform image compression exploiting both spatial and spectral dependencies of hyperspectral image data. The Combined Model allows for the combination of different spectral and spatial autoencoder models. It further allows spatial autoencoder methods that are optimised for image compression on RGB or multi-spectral images to be used in the domain of hyperspectral imaging. We propose multiple spectral and spatial models to be used as part of the Combined Model and a training methodology specific to the Combined Model as well as a novel loss function. The models are evaluated on the large-scale hyperspectral image dataset HySpecNet-11k. We analyse the different aspects of the Combined Model architecture and its training process and compare the achieved reconstruction accuracy for multiple compression ratios. The results show that the Combined Model architecture significantly improves upon state-of-the-art learning-based hyperspectral image compression methods as well as JPEG 2000 for low and very low bitrates.

Zusammenfassung

Kompression von hyperspektralen Bildern mithilfe von neuronalen Netzwerken ist ein wachsendes Forschungsgebiet innerhalb der Fernerkundungsforschung. In dieser Arbeit wird das Combined Model, eine Architektur zur Kompression von Bildern unter Nutzung von sowohl den örtlichen als auch den spektralen Abhängigkeiten in hyperspektralen Bildern, vorgeschlagen. Mithilfe des Combined Models können verschiedene spektrale und örtliche Autoencoder in Kombination miteinander eingesetzt werden. Zudem können örtliche Autoencoder, die ursprünglich für das Gebiet der Kompression von multispektralen oder RGB-Bildern entwickelt wurden, für die Kompression von hyperspektralen Bildern eingesetzt werden. Wir schlagen außerdem mehrere spezifische spektrale sowie örtliche Autoencoder-Modelle vor, die als Teil des Combined Model eingesetzt werden können, sowie eine Trainingsmethodologie, die für das Combined Model entwickelt wurde. Als Teil dieser Trainingsmethodologie wurde zudem eine neue Verlustfunktion entwickelt. Die Modelle wurden anschließend auf dem hyperspektralen Bilddatenset HySpecNet-11k evaluiert. Verschiedene Aspekte des Combined Models sowie die Trainingsmethodologie wurden analysiert und für verschiedene Bitraten verglichen. Die Ergebnisse zeigen, dass das Combined Model für niedrige und sehr niedrige Bitraten signifikant bessere Bildrekonstruktionen erzielt als sowohl bisherige hyperspektrale Kompressionsmethoden auf Basis von Machine Learning als auch die Kompressionsmethode JPEG 2000.

Contents

List of Figures	ix
List of Tables	xi
1. Introduction	1
1.1. Objective	4
1.2. Outline	5
2. Related Work	7
2.1. Hyperspectral Image Compression	7
2.1.1. Traditional Architectures	8
2.1.2. Learning-based Architectures	9
2.2. RGB Image Compression	12
2.2.1. The Hyperprior Architecture	13
3. Theoretical Foundations	15
3.1. Convolutional Neural Networks	15
3.2. Self-Attention modules	16
3.3. Arithmetic Coding	17
3.4. Hyperprior Architecture	19
4. Methodology	23
4.1. The Combined Model	24
4.2. Spectral Autoencoder Methods	26
4.2.1. Per-Pixel Convolutional Neural Network	26
4.2.2. Two-dimensional Spectral CNN	28
4.2.3. Variational Autoencoder	30
4.3. Spatial Autoencoder Methods	31
4.3.1. CNN-based Spatial Autoencoder	31
4.3.2. Hyperprior-based Spatial Autoencoder	32
4.3.3. Attention-based Model Using Hyperprior Architecture	32

5. Experiments	34
5.1. Description of the Data Set	34
5.2. Loss Functions and Metrics	35
5.2.1. MSE Loss	35
5.2.2. Rate Distortion Loss	35
5.2.3. Dual MSE Loss	37
5.2.4. Metrics	37
5.3. Impact of Hughes Phenomenon	38
5.4. Analysis of Dual MSE Loss And Latent Spaces	41
5.5. Analysis of Hyperprior-based Models	44
5.6. Analysis of Kernel Size Change	49
5.7. Analysis of Per-Channel Autoencoder	49
5.8. Rate Distortion Comparison	50
6. Discussion and Conclusion	54
Bibliography	57
A. Appendix	64

Acronyms

1D One-dimensional

2D Two-dimensional

3D Three-dimensional

ANN Artificial neural network

ASCII American Standard Code for Information Interchange

ASI Italian Space Agency

bpppc Bits per pixel per channel

BTC Block Truncation Coding

CABAC Context-adaptive binary arithmetic coding

CAF Cross-layer adaptive fusion

CCSDS Consultative Committee for Space Data Systems

CNN Convolutional neural network

DCT Discrete cosine transform

DFT Discrete Fourier transform

DualMSE Dual Mean Squared Error

DWT Discrete wavelet transform

EnMAP Environmental Mapping and Analysis Program

ESA European Space Agency

GAN Generative adversarial network

GDN Generalized divisive normalization

GMM Gaussian mixture model

GPU Graphical processor unit

GSD Ground sampling distance

H.264/AVC Advanced Video Coding

HEIC High Efficiency File Format

HEVC High Efficiency Video Coding

IGDN Inverse generalized divisive normalization

JPEG Joint Photographic Experts Group

KL Kullback-Leibler

KLT Karhunen-Loeve transform

LeakyReLU Leaky rectified linear unit

MLP Multi-layer perceptron

MSE Mean Squared Error

NLP Natural Language Processing

PCA Principal component analysis

PNG Portable network graphics

PReLU Parameterised Rectified Linear Unit

PRISMA PRecursores IperSpettrale della Missione Applicativa

PSNR Peak signal-to-noise ratio

RD Rate Distortion

RGB Red, green and blue

S2 Sentinel-2

SAM Spectral Angle Mapper

SVM Support vector machine

SWIR Shortwave infrared

TD Tucker decomposition

VAE Variational autoencoder

VNIR Visible and near infrared

List of Figures

3.1. Example of convolutional layer	16
3.2. Arithmetic coding process of the character sequence "ABAC".	18
3.3. Scale hyperprior architecture	20
3.4. Example of masked convolutional layer	21
3.5. Joint hyperprior architecture	22
4.1. Combined model architecture	24
4.2. Per-Pixel spectral convolutional neural network (CNN)	27
5.1. This graph shows the PSNR for the model FastConv and FastCombined (see Table 5.1) for different subsections of the dataset. Only training and validation set sizes were reduced, all models were evaluated on the full test set.	40
5.2. This graph shows the spectral angle for the model FastConv and FastCombined (see Table 5.1) for different subsections of the dataset. Only training and validation set sizes were reduced, all models were evaluated on the full test set.	41
5.3. This graph shows the PSNR for different values of the parameter λ for the Dual Mean Squared Error (DualMSE) loss function.	42
5.4. This graph shows the Inner MSE Loss for different values of the parameter λ for the Dual Mean Squared Error (DualMSE) loss function. Inner MSE Loss is the MSE between the output of the spectral encoder and the reconstruction of that latent from the spatial autoencoder in the Combined Model.	43

5.5.	This image shows four of the 13 channels in the latent representation of an image from the HySpecNet-11k dataset for the model FastCombined-AllFrozen for the values 1 and 0.7 for the parameter λ of the Dual Mean Squared Error (DualMSE) loss function. The input is on the left. Below it is the latent produced by the spectral encoder, which is shared by both models. Following are, in order, the inner latent produced by the spectral encoder for $\lambda = 1$, and for $\lambda = 0.7$, the reconstruction of the latent produced by the spatial decoder for $\lambda = 1$ and for $\lambda = 0.7$. The full image reconstructions are shown on the bottom.	45
5.6.	Inner latent comparison	46
5.7.	Rate distortion plot showing PSNR for different bitrates for multiple model types.	50
A.1.	Dual MSE lambda plot	65
A.2.	Latent image comparison	65

List of Tables

5.1.	Acronyms for commonly used models in this section.	39
5.2.	Comparison of hyperprior-based models. F is the spatial compression/de-compression factor in the main encoder/decoder pair.	46
5.3.	Comparison of different kernel sizes for the spatial autoencoder and per-channel spatial autoencoder with the baseline model. For all models the spectral encoder is frozen after pretraining and $\lambda = 0.5$	48

1. Introduction

Hyperspectral imaging is a quickly growing field. It is the technique of capturing images with a specialized camera in order to obtain a spectrum of many wavelengths of light for each pixel of a taken image. There are three categories of hyperspectral cameras that are used to capture such images. The first, push broom scanners, use a linear arrangement of spectroscopic sensors, that are sensors able to capture a spectrum of many wavelengths of light at once. The sensor arrangement is then moved over the subject of the image, for example by being attached to a satellite orbiting earth. The second category, whisk broom scanners, functions similarly to push broom scanners with the difference being that the stationary sensor arrangement is replaced by a moving mirror reflecting light into a single detector that collects the spectral information in a step-by-step manner. Lastly, snapshot hyperspectral imaging works by using a sensor array to capture a complete hyperspectral image with a single activation of the sensors.

Regardless of capture method, there are many applications for the use of hyperspectral images since a multitude of information can be extracted from their combination of spatial and spectral data. One such application is pixel-wise classification of materials, that being the process of determining the material of specific pixels in a hyperspectral image. This is possible because of the high amount of information present per pixel resulting from the large amount of spectral information. An instance of this approach can be seen in Zea et al. [1]. They use classification of hyperspectral images in order to detect the presence of the toxic metal cadmium in soil, lessening the need for chemical methods that require the plant to be harvested to test for cadmium stress. Using hyperspectral imaging the detection of cadmium can be performed on living plants. Another example of pixel-wise classification materials is explored in Henriksen et.al. [2], where it is used to separate plastic waste by twelve kinds of plastics with a higher accuracy than previously used methods such as near-infrared technology.

Another field where the use of hyperspectral imaging is rising in importance is in geology and environmental sciences. An example of this is the Environmental Mapping and Analysis Program (EnMAP) mission which launched a satellite into earths orbit in 2022. It takes images of the surface of the earth with a comparatively high ground sampling distance (GSD) of 30 square meters per pixel, allowing for regional geographic analysis

[3]. In the spectral domain, the satellite yields a high resolution. By combining a sensor in the visible and near infrared (VNIR) spectrum and a sensor in the shortwave infrared (SWIR) spectrum it can capture light with wavelengths between 420 and 2450 nm. Wavelengths between 420 and 1000 nm are captured by the VNIR sensor with a resolution of 6.5 nm per spectral band, while wavelengths between 900 and 2450 nm are captured by the SWIR sensor with a resolution of 10 nm per band. The EnMAP mission has already lead to many interesting studies including of glacier ice surface properties of the ice sheet in South-West Greenland, moisture content of soil below grassland and identification of specific crop traits such as the chlorophyll content or the leaf water content [4–6]. A dataset created from the images produced by the EnMAP mission also serves as the main dataset studied in this thesis.

Another ongoing hyperspectral imaging mission is the PRISMA (PREcursore IperSpettrale della Missione Applicativa) mission led by the Italian Space Agency (ASI) that was launched in 2019 [7]. Similar to the EnMAP mission, the PRISMA mission also consists of a satellite that observes earth using both a sensor for the VNIR spectrum and for the SWIR spectrum. It also has a GSD of 30 square meters per pixel and a slightly lower spectral resolution given as less than 12 nm per band [8, 9]. In contrast to the EnMAP mission the satellite also carries a panchromatic camera, meaning that the camera only has a single spectral band. This camera however has a higher GSD of 5 square meters per pixel. The data from the PRISMA mission has been used to increase spatial resolution in hyperspectral images. This was done by combining the original hyperspectral images with multispectral images with a higher spatial resolution in a process called hyperspectral-multispectral fusion [10]. The term 'multispectral image' refers to images with more spectral bands than traditional red, green and blue (RGB) images, but less bands than hyperspectral images. Here the multispectral images contain 10 bands and are obtained from the Sentinel-2 (S2) project, a mission by the European Space Agency (ESA) that continually performs multispectral imaging with fine spatial resolution [11]. Another use of the PRISMA data has been in disaster monitoring where it has been used to segment wildfires into areas of active fire, smoke, burned ground, bare soil and remaining vegetation [12]. Another related system developed using PRISMA data performs wildfire detection in real-time from satellites [13].

The increasing usage of hyperspectral imaging makes it essential to address the disadvantages of the technology. One major disadvantage of this type of imaging is the required amount of disk space for the resulting images. Because there are many more spectral bands compared to the three bands of RGB in traditional photography the resulting file size rises accordingly. Furthermore, hyperspectral images often use a high precision for the specific brightness values captured for each point in the image. The EnMAP satellite images are composed of 32 bit values, while in RGB imaging 8 bits per pixel per band is

most common [3]. A hyperspectral image might for example have 300 bands. In combination with 32 bit values per pixel per band the resulting file size would be 400 times larger than an RGB image with the same pixel density. In addition to this it is important to notice that even RGB images need to be compressed for many use cases in order to be efficiently stored and transmitted, which is why there are many compression standards for RGB images in wide usage. Furthermore, compression of hyperspectral images is a complex problem. This arises from the combination of spatial image compression complexities and the added challenge of encoding information in the spectral dimension which RGB compression algorithms do not consider. The latter is one of the reasons why many RGB compression algorithms do not perform well on hyperspectral data as will be shown in this thesis. In addition to this there are more technical reasons for the difficulty of adapting algorithms for RGB compression that will also be explored in this thesis. For the above mentioned reasons research into compression algorithms for hyperspectral images is of vital importance.

Algorithms for image compression can be categorized into multiple broad groups. They can firstly be grouped by whether they compress the image losslessly or with loss. Lossless algorithms restore the compressed image exactly, whereas lossy reconstructions cause distortion. The disadvantage of lossless compression is however that there are mathematical limits to the compression rate given by the entropy in the data that is to be compressed. This limit is given by Shannon's source coding theorem [14], which has been stated as follows [15]:

N independent identically distributed random variables each with entropy $H(X)$ can be compressed into more than $NH(X)$ bits with negligible risk of information loss, as $N \rightarrow \infty$; but conversely, if they are compressed into fewer than $NH(X)$ bits it is virtually certain that information will be lost.

This means that the average bitrate, the compression rate given by the quotient of input bits and output bits of the compression algorithm, which is achievable using lossless compression algorithms is given by the entropy of the data that is to be compressed. In contrast, lossy compression methods can achieve much higher compression rates by allowing the introduction of distortion. Furthermore, they can also adapt their rate of compression based on either the desired amount of distortion or the target compression rate. They can also be combined with a lossless compression method to further optimise their compression rate. Another important factor that determines the trade-off between lossless and lossy compression methods is that many applications do not require the perfect reconstruction given by lossless methods. For example, García-Vílchez et al. [16] showed that in hyperspectral image classification lossy compression of the image does not always reduce the performance of the classifier. For some lossy compression algorithms

the classifier even produced better results on the compressed data than the uncompressed data. This is explained by an introduction of smoothing from the compression which is advantageous for the support vector machine (SVM) classifier used in the experiment. For these reasons this thesis focuses on lossy compression methods.

Lossy compression algorithms can be further categorized into traditional compression methods and learning-based compression methods. Traditional compression methods employ algorithms such as the discrete cosine transform which is used in the Joint Photographic Experts Group (JPEG) image compression standard or the wavelet transform, often in combination with a lossless entropy coding method. Examples of traditional methods other than JPEG are Block Truncation Coding (BTC) which splits images into blocks and then rounds pixel values inside the blocks efficiently and the High Efficiency File Format (HEIC) standard, which uses wavelet transforms similar to JPEG [17, 18]. Learning-based approaches use the powerful capability of artificial neural networks to universally approximate arbitrary functions using gradient descent [19]. These networks are then used to build models that learn to compress and decompress images. Lately these networks have been shown to outperform traditional compression methods for many applications, especially for RGB images using for example the hyperprior model given by Ballé et al. which will be discussed in detail in Section 2.1 [20–22]. Research into compression for hyperspectral images is in its early stages, however even in this domain there are promising results showing the capabilities of learning-based autoencoders for this task [23–26].

1.1. Objective

This thesis addresses the problem of hyperspectral image compression using machine learning models consisting of two parts, an encoder and a decoder. The encoder learns to map the original image to a low-dimensional latent space, thereby performing compression. Analogously, the decoder is trained to reconstruct the input by mapping elements from this latent dimension to full-size images that are as close as possible to the original image.

Contributions to the research on learning-based hyperspectral image compression are made in this thesis by introducing a new architecture that enables the use of spatial compression algorithms such as models that might be used for RGB image compression by combining them with a model that performs compression in the spectral domain while keeping spatial relationships intact. In this way, it is possible to compress the spatial information in hyperspectral images using models that are not ordinarily applicable to

such images.

Using this architecture, multiple model combinations are designed and compared with each other as well as with the base versions of these models. These models include models based on convolutional neural networks (CNNs) as well as transformer-based architectures. A model using a hyperprior architecture is also used for the spatial compression. This architecture yields much higher compression ratios than other models by using an arithmetic encoder in the bottleneck of the compression model. In this way compression rates much higher than the current state of the art for hyperspectral image compression are achieved while distortion is only reduced by a comparatively small amount.

Additionally, the latent spaces of both the encoder in the spectral domain as well as the latent space of the spatial encoder will be analysed.

1.2. Outline

This section gives a brief introduction of the 7 chapters into which this thesis is split.

Chapter 2 presents the related work. This relates to the research of learning-based hyperspectral image compression by itself. However, as learning-based compression for hyperspectral images, in contrast to traditional hyperspectral compression methods, is a recent and relatively unexplored field of study, the adjacent research topic of learning-based RGB image compression is also explored. There is also some exploration of the general studies done on convolutional neural networks as well as transformers.

Chapter 3 gives an overview of the theoretical ideas used in the proposed methods, such as CNNs, transformers, arithmetic coding as well as the hyperprior architecture for compression.

Chapter 4 details the individual submodels making up the models that are proposed to address the hyperspectral image compression problem as well as the models in their totality. It also gives an overview over the loss functions used for training the models, one of which is a loss function specifically designed for the models proposed in this thesis.

Chapter 5 explains the design of the performed experiments, the dataset that is used for these experiments as well as the results of these experiments. It also details some of the

challenges encountered during the experimentation phase.

Chapter 6 gives a summary of the results from the thesis as well as possible improvements that could be made as well as ideas that could be explored in future research.

2. Related Work

Learning-based hyperspectral image compression is a developing field of study. While there are papers published on this topic, there are some problems making it difficult to assess and compare the results of these studies, as will be explored further in this chapter. An overview of hyperspectral image compression algorithms by Dua et.al. [27] shows the discrepancy between traditional transform-based and prediction-based compression techniques and techniques based on machine learning. The review contains 21 transform-based and 19 prediction based techniques but only five learning-based models. While the overview over the models based on machine learning is no longer complete since the review was done in 2020, it still shows that learning-based approaches for hyperspectral image compression until recently only received limited attention from the research community. Lately however, research in this field has been increasing quickly. Since learning-based RGB image compression is a closely related field and much more widely researched, relevant studies done regarding this topic are also analysed.

2.1. Hyperspectral Image Compression

Most learning-based hyperspectral image compression research employs a CNN-based model architecture to reduce the dimensions of the input image [23–25]. Another model proposed by Guo et.al. [26] uses a hyperprior architecture, originally developed by Ballé et.al. [20], which also uses convolutional layers in an artificial neural network (ANN) but combines them with an arithmetic coder to improve compression rate. Some other models that will be explored later are a model based on a SVM by Aidini et.al. [28], a generative adversarial network (GAN) model by Deng et.al. [29] and a model using a simple multi-layer perceptron (MLP) by Kumar et.al. [30]. Before detailing these methods, an overview of some traditional hyperspectral image compression architectures is given.

2.1.1. Traditional Architectures

Dua et. al. [27] describes eight categories of hyperspectral image compression techniques, one of which are learning-based approaches. Of the traditional methods, that is methods not based on machine learning, the most researched categories are transform-based techniques and prediction-based techniques. Transform-based techniques work by transforming the spatial dimension, the spectral dimension or both into the frequency domain by using a transformation function that is applied to the input image [27]. Possible transformation functions are the Karhunen-Loeve transform (KLT), the discrete cosine transform (DCT), the discrete Fourier transform (DFT) or a discrete wavelet transform (DWT). The transformation can then be used to decorrelate the image in the spatial, spectral or both domains depending on the specific technique used. More specifically, the KLT directly decorrelates the data, while for the other transformation functions additional algorithms are commonly used to decorrelate the coefficients resulting from the transformation [31, 32]. After decorrelation, the coefficients are quantized by removing coefficients that are close to zero [27]. The quantized coefficients are then compressed further using an entropy coding algorithm in the final step of compression. Decompression is possible because the used transformation functions are invertible, although loss is introduced during the quantization step. Some of these general steps may be modified for specific transform-based architectures.

Looking at a specific example of a transform algorithm, Karami et.al. [32] use a two-dimensional DWT in combination with the Tucker decomposition (TD) algorithm, an extension of principal component analysis (PCA), in order to reduce correlation. They apply the 2DWT to each spectral band of the input image and obtain four tensors containing the coefficients for each band. The TD algorithm is then used to decorrelate the coefficients both in the spatial and spectral domain. The least significant coefficients are then removed using an iterative process before the remaining coefficients are compressed using the entropy coding technique of arithmetic coding.

The other common category of traditional hyperspectral image compression techniques are prediction-based approaches. These approaches compress images by predicting pixel values from the other values of the other pixels using some algorithm and then only storing the prediction residuals [27, 33]. The current compression standard proposed by the Consultative Committee for Space Data Systems (CCSDS), CCSDS 123.0-B2, is a prediction-based approach [34]. The predictor for this standard predicts a pixel value by computing the sum of its neighbours [33]. The difference between these sums and the original pixel values are then stored in a local difference vector which can then be compressed using an entropy coder.

2.1.2. Learning-based Architectures

In contrast to the traditional methods, learning-based architectures use ANNs trained using gradient descent to learn a mapping from the space of input images to a latent space [19]. In image compression, this latent space is lower-dimensional than the input space, thereby performing compression. The most common approach for learning-based hyperspectral compression are CNN-based architectures. These architectures can be split into two categories, the first being the models using two-dimensional (2D) convolutional layers to learn the spatial dependencies of the hyperspectral images [25]. The second category are models using one-dimensional (1D) convolutional layers to learn the spectral dependencies of the input data [23, 24]. This thesis proposes the first purely CNN-based approach using both the spatial and the spectral dependencies of hyperspectral images for compression. The model proposed by Guo et.al. [26] utilises both spatial and spectral dependencies, it does however rely on a hyperprior architecture and is therefore not a purely CNN-based model. This model as well as the hyperprior architecture will be detailed later in this chapter. As mentioned above, comparing the results from these papers directly is difficult. The reason for this is that the models use different datasets since there is currently no accepted standard dataset for hyperspectral image compression. Furthermore, the models are designed for different compression rates. Since a higher compression rate also leads to a higher distortion for the same model in most cases as described by rate-distortion theory [35], this makes it hard to directly compare the results from papers that use both a different dataset and different compression rates, therefore requiring a reimplementation of the models in order to be able to fairly compare them.

A model exploiting only the spectral dependencies in the input images is provided by Kuester et.al. [23, 24]. They use a model consisting of a 1D convolutional layer to learn spectral dependencies followed by a max pooling layer applied to the spectral dimension to reduce the dimensionality of the image. This block of 1D convolutional layer and a max pooling layer is repeated, using leaky rectified linear unit (LeakyReLU) as the activation function. Following that, two more 1D convolutional layers are added, now without max pooling layers inbetween. Each convolutional layer uses less filters than the layer before it. Finally, the last layer uses only one filter, thereby providing the bottleneck of the compression algorithm with a fixed compression rate of 4. Decompression is performed using a reversed model, symmetrical to the encoder. The decoder only differs by replacing the blocks consisting of a convolutional and a max pooling layer by transposed 1D convolutional layers with a stride of two to increase the spectral dimensionality instead of reducing it [24]. In an earlier variation of the model, upsampling layers are used instead for this purpose [23]. In contrast to most other models, this model is trained and applied not to a whole image at once, but rather each pixel individually. This has the advantage of reducing the complexity of the task the model has to perform. Instead of compressing a

complete hyperspectral image, the model only learns to compress the spectral signatures of arbitrary pixels from these images. The number of training samples also increases dramatically, as each hyperspectral image consists of many pixels, although each training sample contains much less information than it does for a model that trains on complete images at once. A disadvantage of the approach is that modelling spatial dependencies is not possible using a per-pixel training approach. Furthermore, training is much slower when compared to other models. The reasons for this and more detailed comparisons with other model architectures are explored in Chapter 5.

In contrast to the aforementioned model, La Grassa et.al. [25] propose a model that focuses on the spatial dependencies in the image data. This model is trained and applied on complete hyperspectral images in one step. They use a combination of 2D convolutional layers and max pooling layers applied to the spatial dimension of a input image to reduce the spatial dimensionality while simultaneously increasing the number of filters in the convolutional layers to keep the amount of stored information stable. After multiple of these blocks, the output of the last block is flattened into a single dimension, followed by a linear layer that maps that output into a one-dimensional latent space of the chosen size. The decoder is again a symmetrical, reversed version of the encoder, replacing only the blocks of 2D convolutional and max pooling layers by transposed convolutional layers with a stride of 2. This model has the advantage of having the ability to set a precise bit rate by changing the output dimensions of the final linear layer. However, applying the model to hyperspectral images with many bands is difficult. This is because in the proposed model, the first convolutional layer uses 64 filters. Since this layer is applied to the complete input image, an input image with dimensions (C,H,W) where C is the number of spectral bands and H and W are the height and width of the image respectively is transformed into the dimensions $(64,H,W)$. If C is significantly higher than 64, this leads to large information loss in the first layer. This is common in hyperspectral datasets. This can be seen in the HySpecNet-11k dataset [36], the main dataset used in this thesis, which uses 202 channels. Possible solutions to this problem are explored in Chapter 4.

There are also some approaches that do not use a CNN-based architecture. While some of these architectures use models that contain a CNN in addition to other parts, some architectures use different model types completely. Among these are both different types of neural networks as well as learning-based algorithms not using neural networks. One such architecture is proposed by Aidini et.al. [28]. They use quantization to compress the original image, meaning that the precision of the spatial and spectral values of the image are decreased. Then an algorithm tries to recover the original tensor values by trying to reconstruct low-rank tensors as a constrained optimization problem. Afterwards, a spatial super-resolution algorithm using trained dictionary learning is used to increase

the resolution of this image, after which a classifier is trained on these super-resolved images. While this architecture is interesting, it is not directly used in this thesis as its methodology is very different to the neural network based methodologies used in the thesis.

Another category of models use ANNs to determine parameters for lossless compression algorithms. Shen et.al. [37] use a deep belief network to determine the optimal parameters for golomb-rice coding, a lossless coding algorithm that normally assumes a geometric underlying distribution. Using a neural network to determine the parameter removes that necessity. This core strategy is also used by Guo et.al. [26]. They use a hyperprior architecture to compress hyperspectral images. Hyperprior models use an ANN-based model to transform the image data into a latent space that is commonly lower-dimensional than the input image. Then a second ANN is trained on the latent space to determine parameters for an arithmetic coder, a lossless coding algorithm. In both Guo et.al. and the original paper introducing the hyperprior architecture for RGB images, Ballé et.al. [20], both ANNs are CNNs and the latent space is indeed a lower-dimensional space.

Guo et.al. innovates on the original approach in the fact that they assume a student's T distribution instead of a gaussian distribution for the arithmetic coder and in the design of the first CNN. Their version includes both a spatial and a spectral part in the main CNN, making it the only model to learn both spatial and spectral information for hyperspectral image compression. They achieve this by using 2-D convolutional layers for the spatial domain and 3-D convolutional layers with a kernel size of one to include the spectral dependencies. However, a disadvantage of their model is that it is developed only for datasets with a low amount of channels compared to other hyperspectral datasets with the highest having 102 spectral channels. The approach is also not easily adaptable to datasets with a much higher number of channels. This is because the first convolutional layer of the model uses 192 filters and a stride of two, therefore transforming an input image with the dimensions (C, H, W) , where C is the number of spectral channels and H and W are the height and width respectively, into the dimensions $(192, H/2, W/2)$. For one of the two datasets studied in Guo et.al., the KAIST dataset, which contains 31 spectral channels, this is a large increase in information. For the other dataset, ROSIS-Pavia, with 102 to 103 spectral channels depending on the scene, this is already a loss of information, requiring the first layer to already encode some of the spatial information of the image appropriately in order to not introduce large amounts of distortion. For datasets with more channels this problem increases. Furthermore, because of graphical processor unit (GPU) memory limitations, increasing the amount of filters appropriately is not always possible.

Another model using neural networks is proposed by Kumar et.al. [30]. Instead of CNNs

they use a simple multi-layer perceptron as the decoder for the reason that they use this model for real-time onboard image compression which requires a much more simple model for faster execution speed. Another uncommon trait of their architecture is that they do not use a symmetric model, meaning that the encoder and decoder are mirrors of one another. Instead, they only use an ANN for the decoder and use a low-complexity encoder based on matrix multiplication. Hong et.al. [38] propose a novel architecture as well. They use a transformer that works on a spectral embedding by linearly projecting groups of neighbouring spectral bands to an embedding vector. This improves the capabilities of the network since neighbouring bands in hyperspectral images capture detailed changes in the absorption of the underlying material and therefore contain important information, especially for classification tasks. In addition to this they implement cross-layer adaptive fusion (CAF) to improve exchange of information in the transformer section of the model. This means that they use multiple transformer layers and, in addition to the direct connection between adjacent transformer layers, add connections that skip one layer and connect with the layer after using a special CAF module. The transformers can be applied either per-pixel or for small patches. However, their work is not used in the context of image compression but rather image classification and therefore only contains an encoder combined with an MLP head to classify hyperspectral image pixels or patches based on categories such as "Corn", "Grass Pasture" or "Wheat". This means that an application of this architecture to the task of image compression would require substantial additions to the model.

2.2. RGB Image Compression

While the research into hyperspectral image compression strongly increased in recent years, RGB image compression has been a widely researched field for many years. There are many traditional compression algorithms, some of which are installed on every modern operating system and browser. Examples of these are portable network graphics (PNG), a lossless image compression format as well as JPEG, a collection of compression algorithms, the most common of which performs lossy compression of images. An improved version of JPEG, JPEG 2000, also exists, offering increased compression rate for similar distortion values [39]. However, while these algorithms are very popular, learning-based image compression has outperformed methods such as JPEG in both compression ratio and distortion [20, 21]. Furthermore, many of the ideas in learning-based hyperspectral image compression originate from RGB image compression studies and many of the ideas in RGB image compression have not yet been adapted to the hyperspectral realm. For these reasons RGB image compression papers are directly relevant even for a thesis that only concerns itself with the hyperspectral

domain.

2.2.1. The Hyperprior Architecture

One of the most important recent works in RGB image compression was released by Ballé et.al. [21]. The model proposed in this builds on a previous work [20] using a CNN to reduce the dimensionality of the input image, followed by a quantization of the resulting latent and the usage of an arithmetic autoencoder to losslessly compress the quantized latent. The output of the arithmetic autoencoder is then decoded by a CNN that is symmetrical to the encoder CNN, similar to the models for hyperspectral image compression that were discussed previously. This model already outperformed JPEG and the improved JPEG 2000 on the tested data.

The performance of the arithmetic autoencoder depends on the accuracy of the estimated probability distribution that is used to compress the given data. In this area improvements were made in Ballé et.al. [21] by using a smaller, separate CNN that learns to extract parameters for a good probability distribution estimate from the latent resulting from the main CNN. This network also uses an autoencoder structure, meaning that the estimate can also be transmitted in compressed form as side channel information using only a small amount of space. Training a network using gradient descent for this purpose would not ordinarily be possible since the quantized latents have discrete values resulting in zero gradients everywhere. For this reason the quantization is substituted during training by addition of a small amount of uniform noise to dediscretize the latents.

This hyperprior architecture is used in a large portion of modern learning-based RGB image compression models, among these also the hyperspectral image compression model by Guo et.al. [26] which was discussed in Chapter 2.1.

One such example is a paper by Hu et.al. [40] where the model is generalised to include not only two but an adaptable number of CNNs, where each CNN learns the probability distribution of the CNN before it. This leads to slight improvements in the bitrate of the model while not changing the distortion. The distortion remains unchanged since the only loss occurs within the first CNN which compresses the input image. The other CNNs are only used to improve the performance of the lossless arithmetic coders.

The architecture was also further improved in two ways by Minnen et.al. [22]. The first improvement is a generalisation in the structure of the probability distribution from the original scale mixture of gaussians [41] to a gaussian mixture model. Therefore, the second CNN generating the probability distribution parameters has to predict both a scale and a mean instead of only a scale as before. This allows for a better modeling of the true

underlying distribution and the paper shows that the increase in necessary side channel information is lesser than the improvement created by the improvement in probability distribution prediction. The second new idea is the addition of an autoregressive model over the latents of the first, main CNN. This also improves compression performance as the more structure from the latents can be exploited, it does however also increase the computational costs of the network as autoregressive models cannot be trained in parallel.

Another application of the hyperprior architecture is found in Cheng et.al. [42]. They improve on the hyperprior model including a joint autoregressive model given by Minnen et.al. [22] by introducing self-attention modules into the encoder and the decoder CNN. Self-attention modules are an idea taken from the field of Natural Language Processing (NLP). There it is used as an integral part of the novel transformer architecture, which yields state of the art results in machine translation and other NLP fields [43]. In computer vision, it has been used in order to generate attention masks that enable the model to use more bits for the more important parts of an image such as foreground or high-contrast elements and less bits to the less important parts of an image [44–46]. This approach of using self-attention modules is also the one used by Cheng et al. They model the architecture of their self-attention blocks by using a simplified version of the architecture proposed in Liu et al. [44]. The simplification is achieved by removing the non-local blocks, which reduces the training time by four times in their testing while achieving a similar loss [42].

3. Theoretical Foundations

This thesis employs multiple different types of machine learning and compression methods, such as CNNs, self-attention modules, arithmetic coders and the hyperprior architecture. The following sections briefly explain foundations in these topics in order to aid the understanding of the models proposed in this thesis.

3.1. Convolutional Neural Networks

CNNs are a type of ANN, distinguished from other ANNs by the fact that they use convolutional layers as part of their architecture [47]. They were originally designed for the purpose of image classification, which remains one of the main applications for CNNs [48]. However, CNNs are also used in other fields of machine learning such as NLP. The most common type of CNN in the field of computer vision uses 2D convolutional layers in combination with subsampling or pooling layers and activation functions in order to learn features of an input image.

A 2D convolutional layer takes a 2D image-like input and produces a given number of features each represented as a 2D matrix [48]. These features are computed for a given pixel by convolving the neighbourhood of this pixel with a kernel, which is shown in Figure 3.1. These kernels are comprised of the weights of the convolutional layer and are therefore learned using gradient descent. A pooling layer is then used to reduce the spatial dimensionality of the input. Such a combination of a convolutional and a pooling layer can be repeated multiple times, since the output of both the convolutional layer and most pooling layers are again image-like and can be used as input to additional such layers. Between sets of convolutional and pooling layers nonlinear activation functions are applied in order for the network to be able to learn nonlinear features. The combination of multiple sets of convolutional and pooling layers is central to learning-based image compression and is used in a multitude of the proposed models [20, 21, 23, 24, 26].

Compared to other neural network architectures such as MLPs, CNNs have low complex-

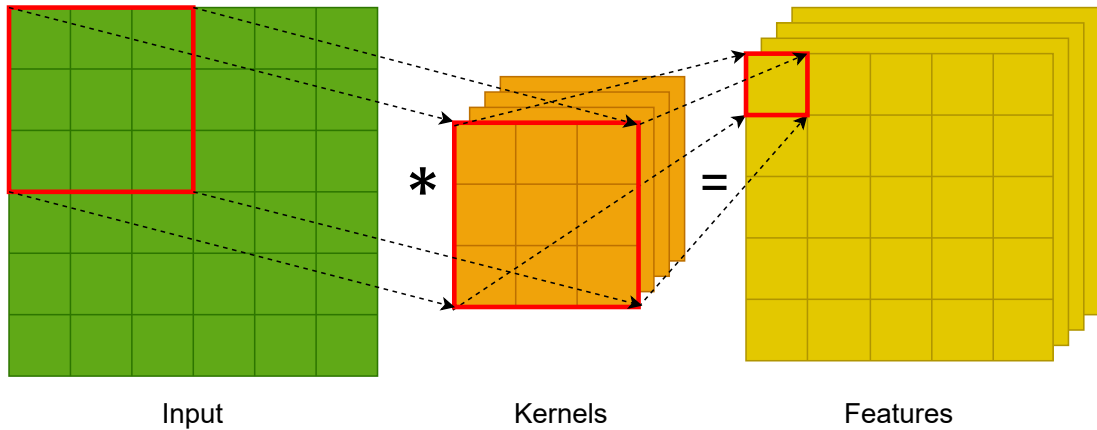


Figure 3.1.: Example of a 2D convolutional layer computing a feature pixel.

ity because of the small number of connections between neurons, as each output neuron of a convolutional layer is only connected to a small number of neurons from the input of the convolutional layer, defined by the kernel size [48]. The low complexity results in CNNs being resilient to overfitting [47, 48]. However, regardless of the low complexity, CNNs are still able to learn high-level features of the input because the use of multiple convolutional layers allows for a hierarchical learning of features. This allows the early convolutional layers to learn to detect low-level features such as corners or edges while the later convolutional layers learn higher-level features [48]

In addition to the 2D convolutional layer, there are other variations of convolutional layers that are useful, in particular for analysing hyperspectral image data. Firstly, one can define a 1D convolutional layer, which differs from the 2D layer in that the kernel moves only in a single dimension. In hyperspectral image analysis this can be used to learn spectral information, since this information is 1D [23, 24]. Similarly, three-dimensional (3D) convolutional layers can be defined by using a 3D cube-shaped kernel. In hyperspectral image processing, this type of layer can be used to learn features of the spatial and spectral domain within a single layer [26].

3.2. Self-Attention modules

The self-attention mechanism is a technique that is widely used in the field of NLP. For example, Vaswani et.al. [43] used it as the main part of their proposed transformer architecture which outperformed other machine-learning approaches in various NLP

tasks at the time, such as language translation. The mechanism aims to compute a latent representation of an input sequence by learning to relate the elements of the input sequence to each other. While a CNN can also achieve this, in a CNN only elements close to one another are used to compute the latent representation. Using self-attention, the relation of all elements with all other elements can be incorporated, enabling learning of long-distance dependencies [43].

In recent years, self-attention modules have also been used in computer vision. In image compression for example, self-attention modules can be used in order to enable the neural network focus on the more complex parts of the input image and subsequently use less bits for the compression of the less complex parts of the input, improving the overall compression[42, 44].

3.3. Arithmetic Coding

Arithmetic coding is an algorithm used for entropy coding, that being a technique for encoding and decoding a sequence of symbols in an efficient way [49]. As an example, a string of characters can be considered. A standard non-optimized method of storing a string is by using the American Standard Code for Information Interchange (ASCII) encoding which used a fixed seven bytes per character of the string. Entropy coding algorithms improve on this by allocating fewer bits to more common characters and therefore more bits to less common characters [49]. In total, this leads to a lower number of required bits for the total message, given an accurate probability model for the frequencies of the characters. In fact, if the provided probability model is perfectly accurate, arithmetic coding yields an encoding close to an optimal encoding, that being an encoding that uses a number of bits equal to the entropy of the encoded input [49]. This is the fewest number of bits possible for a lossless encoding of an arbitrary message, which was proven with Shannon's source coding theorem [14, 15].

Arithmetic coding, using the theoretical version of the algorithm, works by encoding the entire input message into one arbitrary-precision rational number q with $0 \leq q < 1$ [50]. In practice the algorithm is slightly modified to account for the lack of infinite precision in computers, as will be detailed later in this section. The encoding is performed using an iterative process. In order to understand this process, consider for an example an alphabet consisting of three symbols, 'A', 'B' and 'C' with 'A' having the probability 0.5, 'B' having the probability 0.3 and therefore 'C' having the probability 0.2. Encoding a message is performed symbol by symbol. To encode the first symbol, the interval $[0, 1)$ (meaning numbers a with $0 \leq a < 1$) is split into three parts according to these

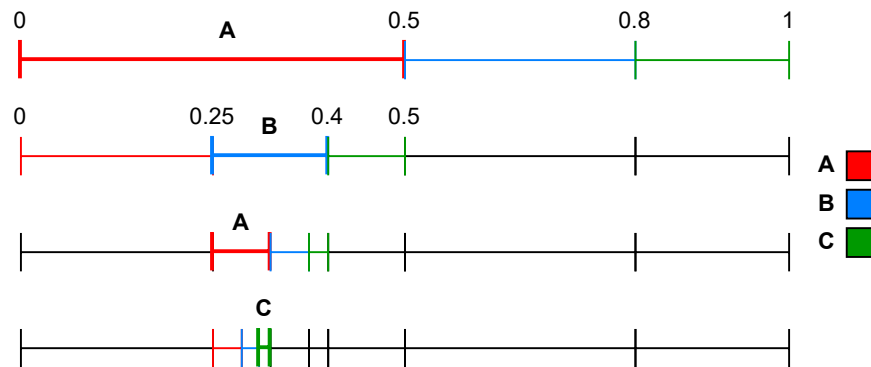


Figure 3.2.: Arithmetic coding process of the character sequence "ABAC".

probabilities. If the first symbol is 'A', the final encoding number q will be in the interval $[0, 0.5)$. Likewise, if the first symbol is 'B' or 'C', q will be in the intervals $[0.5, 0.8)$ or $[0.8, 1)$ respectively. The second symbol is then encoded by splitting the interval determined from the first symbol again, using the same probabilities. A sequence of two 'A' symbols would therefore result in q being in the interval $[0, 0.25)$. This procedure is iterated until each symbol of the input message is encoded, resulting in a single number q representing the entire input. This process can be seen in Figure 3.2 using the aforementioned probabilities for the word "ABAC". This number encodes the original message nearly optimally given a perfect probability distribution, with a difference of less than a bit to the entropy of the message. The reason for this discrepancy is that q is always represented using an integer amount of bits. If the input image has a non-integer entropy, the required amount of bits to store q is the entropy rounded up to the nearest integer. However, the complete algorithm of arithmetic encoding requires an additional source of information that needs to be transmitted to the decoder. In order to decode the message, the decoder needs to know the probability model as well as a way to distinguish when the message ends. The first is constant with respect to message length. The second can be transferred using only logarithmic overhead with respect to message length. This is often done by adding an additional token that signifies the end of the message [50].

As mentioned before, the algorithm is adapted slightly in practice in order to use fixed precision numbers instead of arbitrary precision numbers. Rounding operations as a result of using fixed point numbers are the same for encoder and decoder and therefore do not inhibit the algorithm in principle. However, as a result of rounding, intervals can become small enough that they are not representable with fixed point numbers of the used precision. In this case, the intervals are changed in a determinable way so that they are again representable. This process is called renormalization and multiple algorithms to perform this procedure exist. Arithmetic coding also allows for the possibility of changing the probability distribution during the encoding process based on previously

encoded symbols. This is possible because the decoder decodes the symbols in the same order as the encoder encodes them and can therefore perform the same change to the probability distribution as the encoder, as long as the method for changing the probability model was decided beforehand [50]. An example of such an algorithm is context-adaptive binary arithmetic coding (CABAC), which is an entropy coding method used in the video encoding standards Advanced Video Coding (H.264/AVC) and High Efficiency Video Coding (HEVC) as well as multiple learning-based RGB compression models [20–22]. CABAC is a binary encoding method, therefore requiring the data that is to be encoded to be transformed beforehand. It works by using multiple probability models and, for each element that is to be encoded, selecting one of these models based on previously encoded elements [51]. Since the choice of probability model for each element is only based on elements that were encoded before it, decoding is still possible by performing the same probability model choices as the encoder. Further more, the probability models are updated after each encoding step to adapt them to the data.

In learning-based lossy image compression, arithmetic coding is used as part of the Hyperprior architecture order to encode the latent representation of an input image [20–22]. The latent representation of the input image and the probability distribution used for the arithmetic coder can be learned using ANNs.

3.4. Hyperprior Architecture

The hyperprior architecture was first proposed by Ballé et.al. [21] to improve on an earlier work by Ballé et.al. [20], both in the field of grayscale and RGB image compression. In the earlier work, a CNN using 2D convolutional layers, downsampling layers and generalized divisive normalization (GDN) as the activation function was used to transform an input image into a latent representation. This network is called the analysis transform. The latent representation is then compressed using an arithmetic coder. Afterwards, another CNN is used to transform the decoded output from the arithmetic coder. This CNN, called the synthesis transform, is symmetric to the analysis transform in that it contains the same convolutional layers in inverse order. Additionally, instead of downsampling layers and GDN upsampling layers and the activation function inverse generalized divisive normalization (IGDN) is used. The whole model is then trained using a loss function called rate-distortion loss that optimizes for compression rate and reconstruction distortion simultaneously. Since the arithmetic coder performs lossless compression and the compression rate of the CNN is fixed, optimizing the compression rate corresponds to optimizing the probability distribution estimate of the arithmetic coder and optimizing the distortion corresponds to optimizing the quality of the reconstruction

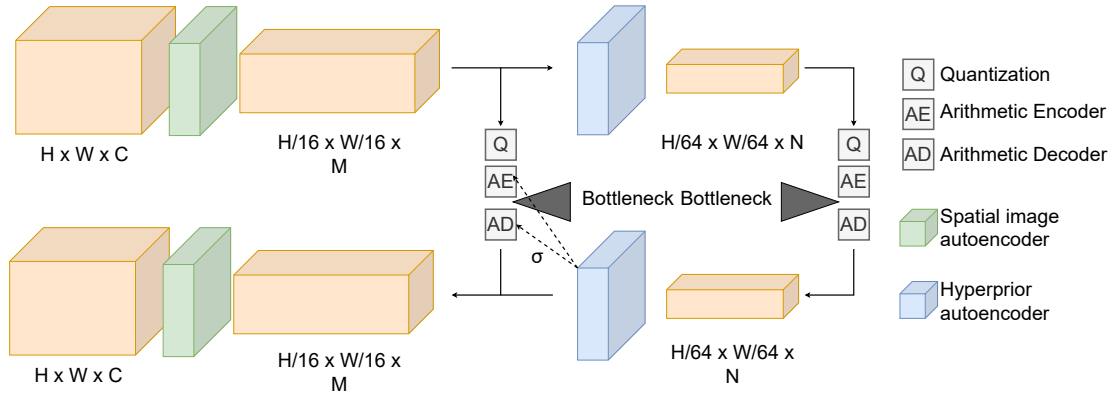


Figure 3.3.: Architecture of the variational scale hyperprior model proposed in [21].

produced by the CNNs. A disadvantage of this approach is that the probability distribution used for the arithmetic coder is a fully factorized distribution, which is not optimal if the latent representation computed by the analysis transform contains statistical dependencies [21].

The hyperprior architecture attempts to solve this problem by transmitting information about the probability distribution of the image that is being compressed as side channel information [21]. This is in contrast to the former model, in which only the probability distribution of the training dataset as a whole is known to both encoder and decoder. The authors noticed that when modelling the probability distribution with gaussians, the standard deviations of latent pixels near each other in specific images tend to vary from the overall probability distribution in a similar manner. Note that while the latent vectors are not images, they can be interpreted as images by treating the filters as different channels of an image with multiple channels. In this way it is also possible to refer to pixels of latents, referring then to the values of all filters for a specific point in the spatial dimensions. To address the spatial dependencies of the standard deviations of the latents, an additional CNN was added to exploit the spatial dependencies in the latents [21]. This network is called "hyperprior". The hyperprior network is also split into an analysis and a synthesis transform and is therefore similar to a smaller version of the original CNN. A visualisation of this architecture is shown in Figure 3.3. A difference between the networks is that the hyperprior network uses LeakyReLU as the activation function instead of GDN. This makes sense because GDN is shown to be especially useful for improving visual similarity in an image compression network. This is not useful for the hyperprior network since it does not attempt to learn image reconstruction. Instead, the hyperprior network is used to learn the probability distribution of the latents of the main CNN. More specifically, the hyperprior analysis transform produces a bottleneck tensor, which can be compressed with an arithmetic coder using a fully factorized distribution

1	1	1	1	1
1	1	1	1	1
1	1	1	1	1
1	1	0	0	0
0	0	0	0	0

Figure 3.4.: Example of a kernel for a 2D masked convolutional layer. The pixels to the bottom and right of the current pixel are not used. The current pixel is marked bold.

similar to the earlier work by Ballé et.al. [20]. The hyperprior synthesis transform then predicts the standard deviations for each pixel in the latent of the hyperprior from that bottleneck. This improves the performance of the arithmetic coder significantly, because the probability distribution predicted by the hyperprior network is closer to the true distribution than the fully factorised distribution that was used in the model without a hyperprior network. This is especially true if the analysed dataset is heterogeneous, since the fully factorised distribution by definition is the same for each compressed image, while the inclusion of a hyperprior enables the arithmetic coder to use different standard deviations for the distributions of each image that depend on that specific image [21]. A tradeoff exists in that in contrast to the non-hyperprior approach, the bottleneck of the hyperprior CNN has to be transmitted alongside the image data. However, Ballé et.al. [21] showed that the advantage of increasing the performance of the main arithmetic coder outweighs the additional bits of information that have to be transmitted.

A further improvement to the hyperprior model was proposed by Minnen et.al. [22], seen in Figure 3.5. They firstly added a context model to the main arithmetic encoder. This context model adjusts the probability distributions for a specific image during the decoding process based on the already decoded elements. This is a common approach in arithmetic coding [50]. See also Section 3.3 for more elaboration on this idea. While the arithmetic coder in the previous models was based on the CABAC algorithm and therefore also context-sensitive, the context model did not improve the model significantly because images that were to be compressed were fed into the autoencoder sequentially pixel by pixel, which is not optimal [20]. Minnen et.al. improves the context model by introducing a masked 2D convolutional layer applied to the quantized latents of the main CNN. A masked convolutional layer in this case refers to a convolutional layer

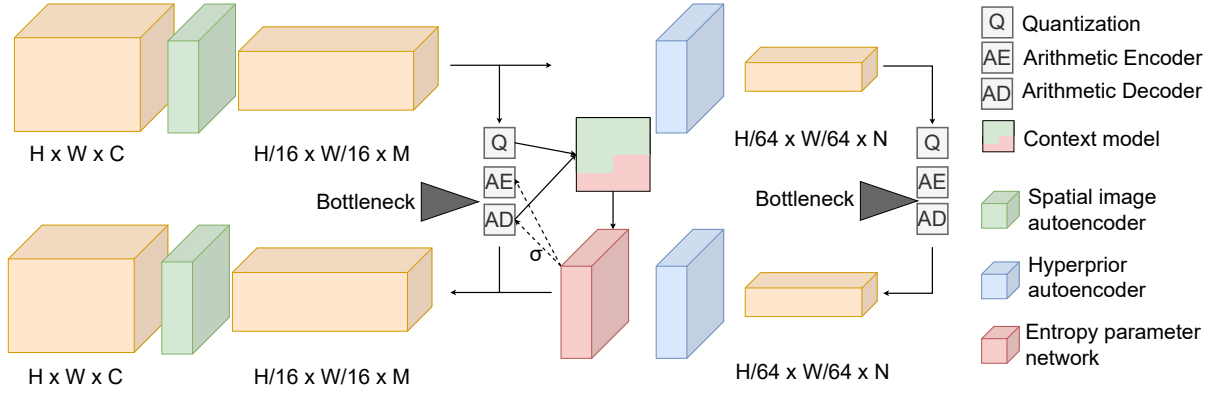


Figure 3.5.: Hyperprior model using an autoregressive context model proposed in [22].

in which the kernel only contains pixels up and to the left of the currently transformed pixel. This idea, first proposed in Oord et.al. [52], ensures the layer only uses the parts of the input that are at this point known to both the arithmetic encoder and decoder. If a standard convolutional layer would be used, the output of the arithmetic coder would not be decodable, since it would be partly based on pixels that are not yet decoded. The output of this layer is then combined with the output of the hyperprior synthesis transform using a three-layer CNN with 1×1 kernels to produce the parameters used for the arithmetic coder. Since the inputs to the context model can also be seen as coming from the output of the arithmetic decoder because of the losslessness of the arithmetic coder, the context model can be seen as an autoregressive component improving the probability distribution estimates of the arithmetic coder [22]. The other improvement proposed in Minnen et.al. [22] is a change from a gaussian scale hyperprior to a gaussian mixture model (GMM). In contrast to the hyperprior in Ballé et.al. [21] that predicted only the standard deviations of the gaussians, the proposed model also predicts the means of the gaussians. This is achieved in the model architecture by making the number of filters in the last layer of the hyperprior synthesis transform and the number of filters in the masked convolutional layer of the context model double the number of filters in the output layer of the main decoder. In this way, for each element in the latent two values in the probability distribution are predicted, the mean and the standard deviation [22].

4. Methodology

As discussed in Chapter 2, learning-based compression of RGB imagery is a hot topic in current research. Therefore, being able to use and further research methods used in that field for hyperspectral image compression would be useful. However, direct application of most learning-based RGB image compression methods on hyperspectral data is not possible. This is mainly for two reasons. Firstly, learning-based RGB image compression methods ignore dependencies between the three channels in the RGB images and focus solely on spatial dependencies. This is useful for RGB images because the spectral dependencies within the three channels are low, since the corresponding wavelengths are far apart. Additionally, having only three channels means that there is not much spectral information per pixel to be compressed, meaning that even if spectral dependencies existed, the compression gains that could be extracted would be small. This is different in hyperspectral image compression where there are large amounts of spectral dependencies and the potential compression gains are much larger because of the high number of spectral channels. It is therefore necessary to modify RGB models to incorporate spectral dependencies in order to adapt them to the task of hyperspectral image compression.

The second reason for the difficulty in applying RGB image compression methods to hyperspectral images is the significant increase in spectral dimension. Compared to an RGB image, a hyperspectral image of the same size with 200 channels has $200/3 \approx 67$ times more dimensions. For this reason RGB models are often too small to learn the compression of high-dimensional hyperspectral image data. For example, many CNN models use an input layer transforming the input from the dimensions $H \times W \times C$, where H and W are the height and width and C is the number of channels to $H' \times W' \times N$, where H' and W' are often either the original height and width or half the original height and width and N is the number of filters used in the input layer. If H' and W' are equal to half the original height and width, this is achieved either using a pooling layer or a convolutional layer with a stride of two. In RGB models, 128 and 192 are common values for N . For RGB images, this means that the input layer increases the dimensionality drastically in order to learn multiple features of the image. However, in hyperspectral images, the same layer would result in a decrease of dimensionality leading to a loss of information. A simple approach to address this issue is to increase the number of filters in each layer according to the increase in dimensionality between RGB and hyperspectral

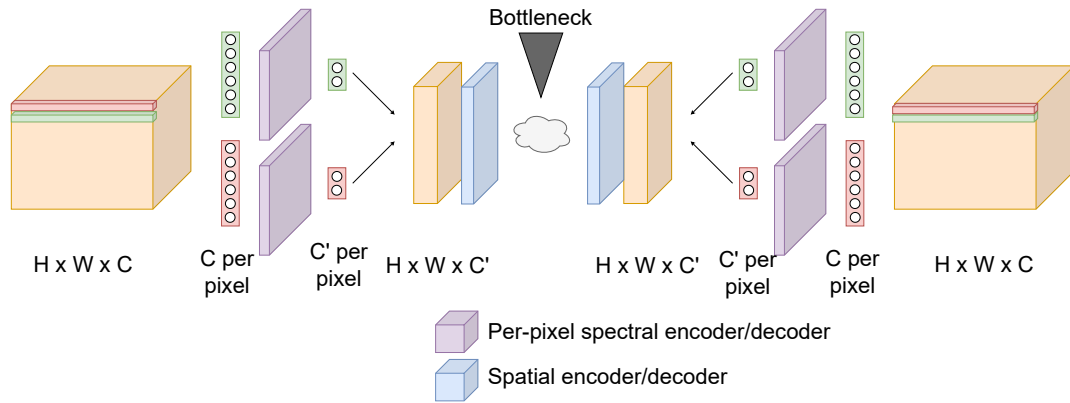


Figure 4.1.: Architecture of the Combined Model using the per-pixel spectral CNN.

images. The input layer would then for example have $192 * 67 = 12864$ filters. However, this would result in a much larger model than the corresponding RGB compression model and therefore require a higher amount of computational resources for training, potentially making the model impossible to train on the available hardware. Additionally, even if the models were trainable, it is unclear if the convolutional layer would be able to learn such a high amount of distinct features for hyperspectral images.

Therefore, a different model architecture is required. This thesis presents a novel architecture that allows the use of spatial compression methods designed for an almost arbitrary number of channels in hyperspectral image compression while also incorporating spectral dependencies. This not only allows the use of RGB image compression methods but also multispectral compression methods. The proposed architecture is called "Combined Model" and an overview is shown in Figure 4.1.

4.1. The Combined Model

The main idea of the combined model is to use a spectral encoder network that preserves spatial dependencies in order to reduce the number of spectral channels in the input. Then, a spatial compression method is applied on the output of that spectral encoder in order to exploit the spatial dependencies. The resulting bottleneck is then decoded using the corresponding spatial decoder. Finally, a spectral decoder network symmetrical to the spectral encoder restores the original image. During training, the spectral decoder learns to adapt to the distortion introduced by the spatial autoencoder. If the spatial autoencoder is learning-based, it similarly learns to produce an output in such a way that

the spectral autoencoder can decode it more effectively. The corresponding mechanism is studied in detail in Chapter 5. Adaptation in either the spatial autoencoder or the spectral autoencoder are necessary, especially if the latent space of the spectral autoencoder is not continuous. In that case, a small change introduced by the spatial autoencoder could result in a large change in the output of the spectral decoder. In practice the adaptation process is successful in avoiding chaotic outputs of the spectral decoder, which is also shown in Chapter 5.

The aforementioned approach has multiple advantages. It is highly flexible with regard to the spatial autoencoder. Many different learning-based approaches can be used and the number of input channels to the spatial autoencoder can be modified so it best fits the spatial autoencoder. This can be achieved by modifying the compression ratio of the spectral encoder/decoder pair, which is possible for all spectral autoencoder models studied in this thesis. Another advantage is that the distinct separation between the spatial and the spectral encoding allows for informative studies to be done both regarding explainability and regarding the amount of spatial and spectral dependencies in the data. The model also allows for different spectral encoder and decoder architectures. There are two necessary conditions required for a spectral encoder/decoder pair to be possible for use in the combined model. Firstly, the output of the encoder must be an image-like shape in order for spatial image compression models to be used for incorporating the spatial dependencies. Secondly, the spectral encoder needs to preserve spatial dependencies, ideally resulting in an output with similar spatial dependencies to the actual input image. Because the spatial autoencoder models were designed for image compression, having spatial dependencies close to the dependencies found in images means that the task the spatial autoencoder is used for is close to the one it was designed for.

However, a difference between the task of image compression and the task the spatial autoencoder performs, that being compression of the latent of the spectral autoencoder, is still present. This is one of the disadvantages of the combined model approach and therefore requires the spatial autoencoder to be resilient to this change of task. However, for the spectral autoencoder models studied in this thesis, the spatial dependencies of the latents very closely resemble the spatial dependencies of the original input, mitigating this disadvantage. For further detail on this see Chapter 5.

For training the Combined Model, an additional idea is useful. If the model is trained in a standard way by training both the spatial autoencoder and the spectral encoder/decoder pair end-to-end from a newly initialized state at the same time, the results can be unstable. Depending on the hyperparameters, it may not be possible to train the combined model in this way. This follows from the fact that inputs to the spatial autoencoder depend on the outputs of the spectral encoder. During the start of training, the spectral encoder

often outputs similar or equal values for all inputs. The spatial autoencoder may then specifically learn to compress latents with these specific values, which is a task very different from compressing spatial dependencies. In this case, the dependencies between the model parts can inhibit training and lead to the model being stuck in the local minimum of outputting the same output for any input. We solve this problem by pretraining the spectral encoder/decoder pair as an autoencoder without any spatial component. Then the combined model is trained using the pretrained spectral encoder and decoder. In this way, the spatial autoencoder receives useful input data from the beginning of training.

An interesting modification option for the model follows from pretraining the spectral encoder/decoder pair. After pretraining the spectral component, the combined model can be trained while freezing the spectral component, training only the spatial autoencoder. If trained in this way, the combined model can be seen as a spatial image compression model using an involved preprocessing and postprocessing step. This option is analysed and compared with the normal method of training the combined model in Chapter 5.

4.2. Spectral Autoencoder Methods

Three different architectures were studied for the spectral encoder and decoder models. Firstly, we analysed a per-pixel CNN using one-dimensional convolutional layers on the spectral domain of the hyperspectral images. We then developed another CNN-based model using two-dimensional convolutional layers to address some of the disadvantages of the one-dimensional CNN approach. Finally, we studied the possibility of directly using a variational autoencoder (VAE) for compression of the spectral information in hyperspectral images.

4.2.1. Per-Pixel Convolutional Neural Network

The first of the spectral autoencoder methods is based on the model proposed by Kuester et.al. [23, 24]. In their proposed method a hyperspectral input is split up into its pixels. The spectrum for each pixel then forms one datapoint on which the model is trained. The encoder of the model being trained is a 1D CNN consisting of two sets of a 1D convolutional layer with a large kernel size of 11 and a max pooling layer followed by two additional 1D convolutional layers with a kernel size of 9 and 7. The number of filters is also reduced by a factor of two in each convolutional layer and set to one in

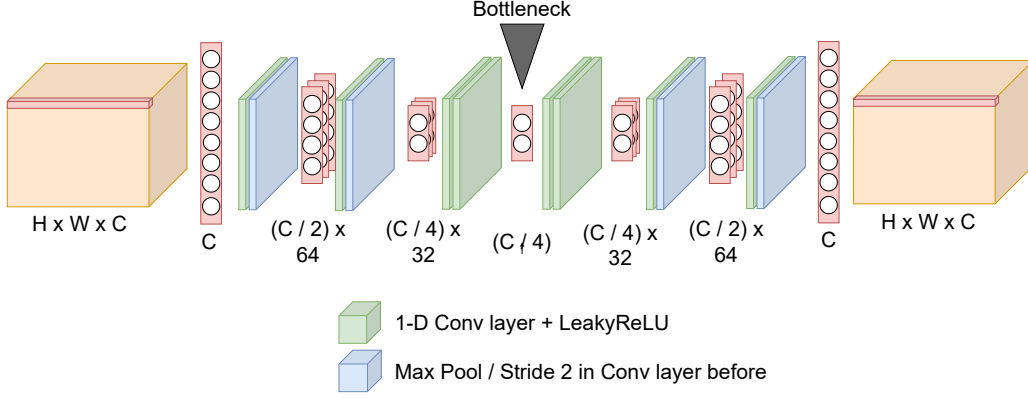


Figure 4.2.: Architecture of the per-pixel spectral CNN using two downsampling layers resulting in a compression factor of four.

the last layer, resulting in a bottleneck with one fourth of the original spectral channels. The decoder is built analogously, only replacing the max pooling layers by upsampling layers. After each convolutional layer a LeakyReLU is used as the activation function. Additionally, the number of channels are padded with channels containing zero values in order for them to be divisible by the downsampling factor of 2. Since the number of channels is reduced by a factor of twice the downsampling factor during encoding, this can result in the bottleneck not having exactly one fourth of the channels of the input, if the number of channels in the padded input is not divisible by four. The decoder then outputs one channel more than the original input. In this case, the last channel of the output is removed. A whole image can be compressed by encoding and decoding each pixel of the image separately, then reassembling the output image from the decoded pixels. The architecture of this model is shown in Figure 4.2.

This model is highly compatible with the Combined Model approach. Because the encoding of the spectral channels is performed per pixel, spatial dependencies are guaranteed to remain during encoding. The latent image is produced similarly to the output of the decoder by reassembling the latents for each pixel into a latent image, which can then be further encoded using a spatial autoencoder. In order to optimize the model for use in the combined model architecture and for the used dataset HySpecNet-11k [36], we adapted the model in various ways.

In the model as proposed by Kuester et.al. [23, 24], the number of layers and therefore the compression ratio is fixed to four. However, different datasets can have different amounts of spatial and spectral dependencies. For of the Combined Model, this can result in the optimal model requiring a different spectral compression factor than 4. To allow for a different compression ratio, the number of sets of a convolutional and a max

pooling layer was enabled to be variable. In addition to this, the number of filters in each convolutional layer was adapted to the number of channels in the data. In the original paper, the first layer has a number of filters equal to half the number of padded channels of the data. We adapted the number of filters to the higher number of channels in the HySpecNet-11k dataset by keeping this ratio equal. We thereby ensure that the number of filters is always appropriately divisible regardless of the number of layers by performing the padding on the spectral channels of the input. We also changed the padding strategy. Instead of padding to a number of spectral channels divisible by half the compression ratio and then conditionally removing some number of channels at the output, we directly padded to a number of spectral channels divisible by the compression ratio. Therefore the need for removing channels at the output is removed, simplifying the model without loss of performance.

The spectral compression model is not only interesting for its applicability to the combined model. It also achieves compression with low distortion especially for higher bitrates. For more on this see Chapter 5. However, there are also some disadvantages of the model architecture. The reduction of spectral channels using max pooling layers results in divisions by two in the number of channels, forcing the compression ratio to be a power of two. Additionally, the higher the chosen compression ratio is, the higher the number of needed padding channels is. As an example, for the HySpecNet-11k dataset with 202 channels two padding channels are needed for the compression ratio 4, for a compression ratio of 32 14 padding channels would be required. Finally, the training process of this model is much slower than the training process for models which take complete images as input. Since each pixel is compressed separately, the number of optimization steps per epoch is much higher. This is also further discussed in Chapter 5.

4.2.2. Two-dimensional Spectral CNN

To address the issues of the 1D spectral CNN, we propose another model that uses 2D convolutional layers with a kernel size of one to compress a whole image at once. For convolutional layers with a kernel size of one the convolution operation collapses down to a scalar multiplication. Additionally, the output has the same height and width as the input, as each pixel in the output is calculated only from the channels of a corresponding input pixel. The output for a specific pixel and filter is therefore equivalent to a weighted sum over all channels in the input for this pixel. Importantly however, the weight used for a specific input channel is the same for each pixel in the image. A two-dimensional convolutional layer with a kernel size of one can therefore only encode spectral information when applied to a hyperspectral image, making it usable as part of the Combined Model. Such layers are commonly used for decreasing the

number of filters between other convolutional layers TODO REF. Because of this, they are also referred to as channel or feature map pooling layers. They are then able to learn non-trivial summarizations of the preceding layer because of the nonlinearity introduced by the used activation function, in our case LeakyReLU. However, they can also instead be used to achieve the opposite effect of increasing the number of filters while retaining the semantic information within these filters. An example of both of these uses is present in the "Inception" architecture proposed in Szegedy et.al. [53].

We show that it is also possible to use such layers to learn spectral dependencies in hyperspectral image data. The model is comprised of two-dimensional convolutional layers starting with a number of filters slightly higher than the number of channels in the hyperspectral dataset the model is trained on, which is 256. In each layer, the number of filters is doubled in order to increase the amount of information available. Finally, a last convolutional layer is used to reduce the number of channels to result in the desired bottleneck size. All convolutional layers use a kernel size of one and the LeakyReLU activation function.

The 2D architecture has multiple advantages compared with the one-dimensional spectral CNN. Instead of being applied per pixel, it can be applied to the entire image at the same time. It also requires less parameters, as each layer has a number of parameters equal to $N_p = C * F$, where C is the number of channels of the input to the layer and F is the number of filters in the layer. Both of these factors result in a much shorter training time as seen in Chapter 5. The model is also more flexible with regard to the compression ratio. In contrast to the one-dimensional spectral encoder, where the compression ratios are necessarily given by $R \approx 2^N$, where N is the number of pooling layers, the two-dimensional spectral encoder allows setting of the number of output channels to an arbitrary number. This also removes the need for padding, which is an advantage especially for high compression ratios which would otherwise require padding. Furthermore, the one-dimensional spectral CNN gets slower for higher compression ratios as it requires more layers while the number of layers and therefore the speed of the two-dimensional spectral CNN is constant with respect to the compression ratio.

Two modifications were also tested for this model. First, a Layer Normalization was performed after each convolutional layer in order to improve generalization performance and further reduce training time. We also modified the model by changing the kernel sizes of the convolutional layers. With higher kernel sizes, the model becomes able to incorporate some amount of local spatial dependencies in addition to the spectral dependencies. We used this to study how sensitive the combined model is to spatial dependencies being partially encoded in the spectral encoder/decoder.

4.2.3. Variational Autoencoder

As mentioned in Section 4.1, a requirement for spectral encoder / decoder pair is that the decoder is able to decode the latent of the encoder after some distortion has been to it applied by the spatial autoencoder. Therefore, using a VAE for the spectral compression is potentially advantageous, because VAEs produce a continuous latent space as shown by TODO REF. A VAE learns a probability distribution for a latent space and produces reconstructions by generating a samples from this distribution. We designed a one-dimensional VAE with the purpose of learning to autoencode spectral dependencies for single pixels, similar to the per-pixel CNN. This per-pixel VAE in fact uses the per-pixel CNN to produce a latent for a given input pixel. Two independent linear layers are then applied to this latent to produce a mean and a variance tensor respectively. The dimensionality of these vectors is variable and therefore a hyperparameter of the model. The linear layer producing the mean uses a tanh activation function to constrain the means between -1 and 1. The linear layer producing the variance also uses tanh but follows by an application of the softplus function to produce a positive variance. A new latent is then sampled from this distribution to produce the input for the decoder which is a symmetric inverted version of the encoder with only one linear layer instead of two, since only the new latent needs to be transformed to the latent space of the per-pixel CNN.

A basic VAE has not previously been used for hyperspectral image compression. Especially compression of spectral signatures using a one-dimensional VAE has not been studied. The only VAEs architecture commonly used for image compression is the hyperprior architecture, which is described in Chapter 4. However, these models are not applicable to be a spectral autoencoder for the Combined Model for two reasons. First, hyperprior models so far have been used for spatial dependencies in images, being applied to whole images at once and not one-dimensionally to spectral signatures. Secondly, since hyperprior models use an arithmetic coder in the bottleneck, they can only be used as the last stage of a compression model since the output of the arithmetic coder cannot be further compressed using another model. This model attempts to overcome these limitations to produce a continuous latent space that can be further compressed using a spatial autoencoder method while being permissible with regard to distortions being introduced during this process.

4.3. Spatial Autoencoder Methods

The second part of the Combined Model is the spatial autoencoder that further compresses the bottleneck of the spectral encoder. While there are some preconditions required for the spectral encoder/decoder pair, the choice of the spatial autoencoder is highly flexible. Since the latents of the spectral encoder have similar spatial dependencies to actual images, any autoencoder that is able to compress images can be used for this part of the model. The models analysed in this thesis are a CNN-based autoencoder, a hyperprior model and an optimized hyperprior model using self-attention modules and a context model. We also analysed the use of JPEG 2000 as the spatial autoencoder to study the viability of using a spatial autoencoder as part of the combined model that is not learning-based.

4.3.1. CNN-based Spatial Autoencoder

The CNN-based spatial autoencoder model uses a low-complexity approach using only two-dimensional convolutional layers. A similar architecture was used for hyperspectral image compression by La Grassa et.al. [25]. However, there are multiple differences between the model. La Grassa et.al. use a linear layer as the final layer of their network. This was not done for the CNN-based spatial autoencoder as it would increase the number of parameters of the model and therefore increase both training time and the required computational resources. The model proposed in this thesis also adds regularization into the model using layer normalization [54]. The number of filters in each layer is also different. Further, La Grassa et.al. use max pooling layers to reduce the spatial dimensionality of the latent while we use convolutional layers with a stride. Additionally, we experiment with different kernel sizes to vary the amount of spatial dependencies introduced to the autoencoder.

The CNN-based spatial autoencoder model first uses two-dimensional convolutional layers to increase the dimensionality of an input image. This is done by using a number of filters much higher than the number of channels in the input. This input layer is followed by another convolutional layer with the same number of filters to increase the depth of the model. Following that is a convolutional layer with a stride of two to reduce the spatial dimensions. This layer also increases the number of filters in order to reduce the information loss while downsampling. Finally, two further convolutional layers reduce the filter dimension to be identical to the number of channels in the input. In this way the dimensionality of the input is reduced by a factor of two in each spatial dimension, resulting in an overall compression factor of four. Additional layers with a

stride of two can be added to increase the compression factor. The decoder is symmetrical to the encoder but replaces convolutional layers with transposed convolutional layers. The model uses the Parameterised Rectified Linear Unit (PReLU) activation function after each convolutional layer and the sigmoid activation function after the final layer of the encoder and the decoder to keep outputs in the interval $[0, 1]$. Layer normalization [54] is also performed after each convolutional layer to increase robustness and rate of convergence during training. Layer normalization was used instead of the more common batch normalization because it is able to perform well even with a low batch size and the combined model with the per-pixel CNN was trained using a batch size of one.

This model was deliberately designed to be low-complexity to reduce the confounding factors for studying properties of the combined model. As there is no prior research for combining models in the way it was performed in this thesis for hyperspectral or RGB image compression, the properties of the combined architecture as a whole was a focus of the experiments. A low-complexity spatial model is therefore a useful baseline that allows for better attribution of specific experimental results to the combined architecture.

4.3.2. Hyperprior-based Spatial Autoencoder

To also include a more complex model with high relevancy in RGB compression research, we implemented the scale hyperprior model proposed by Ballé et.al. [20]. Since this model was explored in detail in Chapter 4, we will not go into further detail here. The only necessary adaptation to the model in order to use it as part of the Combined Model was a change of the input and output layer to accommodate the chosen number of channels in the latent of the spectral encoder. This model was not studied intensively but rather used as a baseline. The main hyperprior-based architecture used was the following model using self-attention modules and a context model.

4.3.3. Attention-based Model Using Hyperprior Architecture

The main hyperprior-based model that was studied is based on work by Cheng et.al. [42]. They further improve the joint autoregressive and hierarchical hyperprior model by Minnen et.al. [22] discussed in detail in Chapter 4 by introducing self-attention modules. The attention modules are based on work by Liu et.al. [44]. However, the non-local blocks have been removed in order to simplify the attention modules. This allowed them to reduce the complexity of the model while still being able to adequately capture

long-range dependencies in the data. These self-attention modules were inserted into the main encoder and decoder in two places as shown in Fig. TODO.

We adapted this model in multiple ways for use in the combined model. First, the model was adapted to be able to receive different numbers of input channels by changing the input and output layers. The second modification regards the strength of the spatial downsampling in the convolutional parts of the hyperprior model. In the original model the main encoder reduces the spatial dimensionality by a factor of 8 in each direction, 64 in total, to produce the latent that is input to both the arithmetic coder and the hyperprior encoder. The hyperprior encoder reduces the spatial dimensions by a factor of 4 in each direction, 16 total, again. Since they used a dataset with 256x256 pixels, this results in a bottleneck of size 32x32 after the main encoder and 8x8 after the hyperprior encoder. However, for hyperspectral data in general and the HySpecNet-11k dataset especially a reduction in spatial dimension with a factor of 64 can result in large distortion, since there are less spatial dependencies than in typical RGB datasets. We therefore modified both the numbers of convolutional layers and the strides in the main encoder/decoder pair and the hyperprior encoder/decoder pair in order to allow for different options for the spatial dimensionality reduction factors. For this, we group models by the spatial reduction factors in the main encoder/decoder pair and the bottleneck size of the hyperprior encoder. The bottleneck size of the hyperprior being larger results in more information remaining present in the hyperprior which increases performance. However, more side information needs to be transmitted, which can decrease performance. This is because the loss used for training these models, Rate Distortion (RD) loss, tries to achieve a balance between bitrate and distortion. Since a larger amount of transmitted side information raises bitrate, RD loss with the same parameters will also result in higher distortion. Therefore, an optimum needs to be found depending on the amount of spatial dependencies present in the dataset.

5. Experiments

The chapter provides detailed analysis of the models proposed in Chapter 4. It also introduces a novel loss function, Dual Mean Squared Error (DualMSE) loss, as well as a training procedure optimized for the Combined Model. All learning-based compression methods explored in this thesis use a set of hyper-parameters to determine different aspects of the model and its training process. Most of these hyper-parameters are specific to a model or an experiment and are therefore noted in the experiments they employed in. However, some hyper-parameters are equal for all experiments. In all experiments the chosen optimizer is the Adam optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The learning rate is set to $0.5 \cdot 10^{-5}$ for most experiments. If a different learning rate is used for an experiment, it is noted explicitly.

5.1. Description of the Data Set

The dataset used in this thesis is the HySpecNet-11k dataset proposed by Fuchs et.al. [36]. HySpecNet-11k is a hyperspectral benchmark dataset that consists of 11,483 image patches containing large-scale views of the surface of the earth. The image patches do not overlap which is advantageous for use in learning-based approaches as there are no duplicated image parts in the training set. The source data for the dataset uses 224 spectral bands and 128×128 pixels spatially with a GSD of 30 meters. However, we use the preprocessed version of the dataset which removes 20 spectral bands due to high water vapor absorption in those bands. We therefore only use 202 spectral bands. The image patches are of high quality because the tiles they were obtained from contain less than 10% snow and cloud cover and were radiometrically, geometrically and atmospherically corrected [36]. Furthermore, cropped patches at the borders of tiles were discarded, so each image patch contains information in all 128×128 pixels. For splitting the dataset into training, validation and test set we used the provided patch-wise splitting set. In this splitting set patches from the same tile can be present in different sets, unlike in the tile-wise splitting set. We did not perform experiments with the tile-wise splitting set as it is outside of the scope of this thesis. It is however a potential future research direction.

5.2. Loss Functions and Metrics

5.2.1. MSE Loss

The Mean Squared Error (MSE) function is a common loss function used for training ANNs. For the task of compressing a hyperspectral image it can be defined as follows:

$$\text{MSE}(X, \hat{X}) = \frac{1}{B * C * H * W} \sum_{b=1}^B \sum_{c=1}^C \sum_{h=1}^H \sum_{w=1}^W (x_{bchw} - \hat{x}_{bchw})^2 \quad (5.1)$$

Here X is defined as a batch of hyperspectral input images, \hat{X} the reconstruction of the images in X obtained from the compression model, B the number of images per batch and C, H and W the number of spectral channels, height and width of the hyperspectral images respectively. We use MSE loss for pre-training the spectral autoencoder models as well as for training the CNN-based spatial autoencoder outside of the combined model.

5.2.2. Rate Distortion Loss

In contrast to the purely CNN-based models with a fixed compression rate for a given set of hyper-parameters, the hyperprior-based models have a variable compression rate that changes during the training process. This comes from the fact that the bitrate of the arithmetic coder depends on the performance of the hyperprior CNN. Because of this MSE loss is not optimal for these models. While it is possible to train hyperprior-based models using MSE loss because the bitrate dependency is contained in the lossless arithmetic coder, this leads to very high bitrates. However, using MSE loss for a hyperprior model can still be useful to obtain an upper bound on the reconstruction accuracy for a given set of hyper-parameters. For more elaboration on this see TODO.

To train a model using an arithmetic coder in the bottleneck with both a low distortion and a low bitrate, a different loss function is needed. The type of problem that arises is a rate-distortion optimization problem [21]. This can problem can parametrized in the following way [22]:

$$R + \varphi \cdot D = \underbrace{\mathbf{E}_{x \sim p_x}[-\log_2 p_{\hat{y}}(\lfloor f(x) \rfloor)]}_{\text{rate}} + \underbrace{\varphi \cdot \mathbf{E}_{x \sim p_x}[d(x, g(\lfloor f(x) \rfloor))]}_{\text{distortion}}. \quad (5.2)$$

We parametrize the rate-distortion trade-off with the Lagrange multiplier φ . (We use φ instead of the commonly used λ to avoid confusion with the λ of the DualMSE loss

described in Section 5.4). p_x is the unknown distribution of input images, $\lfloor \cdot \rfloor$ refers to quantization and f is the main encoder. Accordingly, $y = f(x)$ is the output of the main encoder applied to the input x , $\hat{y} = \lfloor y \rfloor$ are the quantized latents of the main encoder, $p_{\hat{y}}$ is a discrete entropy model and $\hat{x} = g(\hat{y})$ is the reconstruction of x obtained from applying the main decoder g to the quantized latents \hat{y} . As the metric for computing the distortion we use MSE, because Ballé et.al. showed that the rate-distortion problem is equivalent to minimizing the Kullback-Leibler (KL) divergence over p_x [21]. Using this equivalence they showed that using MSE as the distortion metric is equivalent to assuming a Gaussian distribution for $p_{x|\hat{y}}(x|\hat{y})$. While there are approaches assuming different probability distributions like the Student's T distribution, this was not a focus of this thesis and the Gaussian often provides a good estimate for unknown probability distributions. For this reason we also assumed a Gaussian distribution for the probability distribution.

However, this formulation does not include the hyperprior network that learns and transmits side information about the entropy model $p_{\hat{y}}$. In order to incorporate the hyperprior into the RD loss function, we need to add an additional term to capture the bitrate for the transmitted side information. The final loss function is then defined as follows:

$$\begin{aligned}
 R + \varphi \cdot D &= \underbrace{\mathbf{E}_{x \sim p_x}[-\log_2 p_{\hat{y}}(\lfloor f(x) \rfloor)]}_{\text{rate (main latents)}} + \underbrace{\mathbf{E}_{x \sim p_x}[-\log_2 p_{\hat{z}}(\lfloor f_h(f(x)) \rfloor)]}_{\text{rate (hyper-latents)}} \\
 &\quad + \underbrace{\varphi \cdot \mathbf{E}_{x \sim p_x}[\|x - g(\lfloor f(x) \rfloor)\|_2^2]}_{\text{distortion}}. \\
 &= \underbrace{\mathbf{E}_{x \sim p_x}[-\log_2 p_{\hat{y}}(\hat{y})]}_{\text{rate (main latents)}} + \underbrace{\mathbf{E}_{x \sim p_x}[-\log_2 p_{\hat{z}}(\hat{z})]}_{\text{rate (hyper-latents)}} + \underbrace{\varphi \cdot \mathbf{E}_{x \sim p_x}[\|x - \hat{x}\|_2^2]}_{\text{distortion}} \quad (5.3)
 \end{aligned}$$

Here z and \hat{z} refer to the latent of the hyperprior and the quantized latent of the hyperprior respectively, that being the result of applying both the main encoder f and the hyperprior encoder f_h to the input x .

It is to note, that since only the distortion metric is dependent on φ , RD loss converges to that metric for $\varphi \rightarrow \infty$, if we assume an appropriately scaled learning rate. Since we use MSE loss for the distortion metric, for our experiments RD loss converges to MSE loss for $\varphi \rightarrow \infty$.

With regard to the quantization operation, care has to be taken. When executing the trained hyperprior model, quantization is defined as rounding to the nearest integer. However, during training this would result in zero gradients everywhere. To avoid this, while training we instead apply a uniform noise to the input in place of rounding. From

this follows that the arithmetic coder is only executed during testing, not during training. In training the achieved bitrate is estimated using a heuristic. However, we verified that the bitrate estimates are close to the real achieved bitrate in Section 5.8.

5.2.3. Dual MSE Loss

The combined model consists of an outer autoencoder model and an inner autoencoder model. However, applying a loss function such as MSE loss in the standard way results in the inner autoencoder being trained through a layer of distortion at both the input and the output, as the input to the loss function are the input to outer encoder x and the output of the outer decoder \hat{x} . The input to the inner encoder, $y = f_o(x)$, where f_o is the outer encoder and the output of the inner decoder, $\hat{y} = g_i(f_i(y))$, where f_i and g_i are the inner encoder and decoder respectively, are not direct inputs to the loss functions. To alleviate that, we designed an experimental loss function called DualMSE loss. It is defined as follows:

$$\text{DualMSE}(x, \hat{x}, y, \hat{y}) = \lambda * \text{MSE}(x, \hat{x}) + (1 - \lambda) * \text{MSE}(y, \hat{y}). \quad (5.4)$$

The DualMSE loss is therefore a linear interpolation between the MSE loss of the complete model and the MSE loss of the inner model. This allows for fine-grained control over the weight of the loss of the inner model in the overall loss. Of note are also the two edge cases of $\lambda = 1$, where the DualMSE loss function becomes equivalent to MSE loss and $\lambda = 0$, where only the loss of the inner model is used. The DualMSE loss was used for the combined models using the CNN-based spatial autoencoder. Results regarding the use of DualMSE loss for these models can be seen in Section 5.4. While the use of a Dual Rate Distortion loss would be possible for the combined models using a hyperprior-based spatial autoencoder, this was not done because the results in Section 5.4 did not show a clear and significant improvement when using a dual loss function.

5.2.4. Metrics

We use two different metrics to evaluate our experiments. The first metric is peak signal-to-noise ratio (PSNR), which is a common metric to quantify reconstruction accuracy for lossy learning-based image compression [20–25]. PSNR is defined in general as follows:

$$\text{PSNR}(x, \hat{x}) = 10 * \log_{10} \frac{\text{MAX}}{\text{MSE}(x, \hat{x})} = 20 * \log_{10}(\text{MAX}) - 10 * \log_{10}(\text{MSE}(x, \hat{x})). \quad (5.5)$$

Here MAX is the maximum possible pixel value in the used dataset and $\text{MSE}(x, \hat{x})$ refers to the Mean Squared Error between x and \hat{x} . However, since the data range for the HySpecNet-11k dataset used in this thesis is the intervall $[0, 1]$, the equation simplifies to:

$$\text{PSNR}(x, \hat{x}) = -10 * \log_{10}(\text{MSE}(x, \hat{x})). \quad (5.6)$$

It is therefore directly correlated with the log of the MSE. While PSNR is a commonly used metric in lossy learning-based hyperspectral compression, it is computed as an average over all channels. We therefore want to employ another metric to more accurately evaluate the reconstruction of the spectral signatures which are important for lossy reconstructions of hyperspectral images.

For this we use the spectral angle. It is computed by calculating the angle between the spectra of the input image and the reconstruction for each pixel using the Spectral Angle Mapper (SAM) algorithm [55]. The computed angles are then averaged to compute the metric for a given input image and reconstruction. Using the combination of PSNR and spectral angle we can compare different models with regard to spatial and spectral reconstruction accuracy.

5.3. Impact of Hughes Phenomenon

Hughes Phenomenon, also referred to as curse of dimensionality, describes the phenomenon that the number of required training samples for a learning-based algorithm is correlated with the dimensionality of the data [56]. This is especially applicable to the domain of hyperspectral imaging because of the high number of spectral channels increasing the dimensionality of the input. We tested the impact of the Hughes Phenomenon on the **FastConv** model and the **FastCombined** model by varying the size of the training dataset. More details on the model labels is provided in Table 5.1. The **FastConv** models were trained for 70 epochs, which resulted in convergence. The trained **FastConv** models were then used as pretrained parts of the **FastCombined** models. The models were then trained for an additional 50 epochs using $\lambda = 1$ while freezing the spectral encoder. The **FastConv** models used a compression ratio of 16 or a bitrate of 2 bits per pixel per channel (bpppc). The **FastCombined** models used an additional spatial compression factor of 4, resulting in a total compression ratio of $16 \cdot 4 = 64$ or a bitrate of 0.5 bpppc.

As seen in Figure 5.1 and Figure 5.2, the models perform worse for smaller dataset sizes. However, for division factors $f < 8$ the gradient decreases with lower division

Table 5.1.: Acronyms for commonly used models in this section.

Model Name	Description
Conv1D	Per-pixel spectral autoencoder (Subsection 4.2.1)
FastConv1D	2D spectral autoencoder (Subsection 4.2.2)
Conv2D	2D CNN-based spatial autoencoder (Subsection 4.3.1)
AttentionHyperprior	Attention-based spatial model using hyperprior architecture (Subsection 4.3.3)
Combined	Combined Model using Conv1D as the spectral part and Conv2D as spatial part (Section 4.1)
FastCombined	Combined Model using FastConv1D as the spectral part and Conv2D as spatial part (Section 4.1)
CombinedWithAttention	Combined Model using Conv1D as the spectral part and AttentionHyperprior as spatial part (Section 4.1)
Spectral-VAE	Spectral VAE-based autoencoder (Subsection 4.2.3)

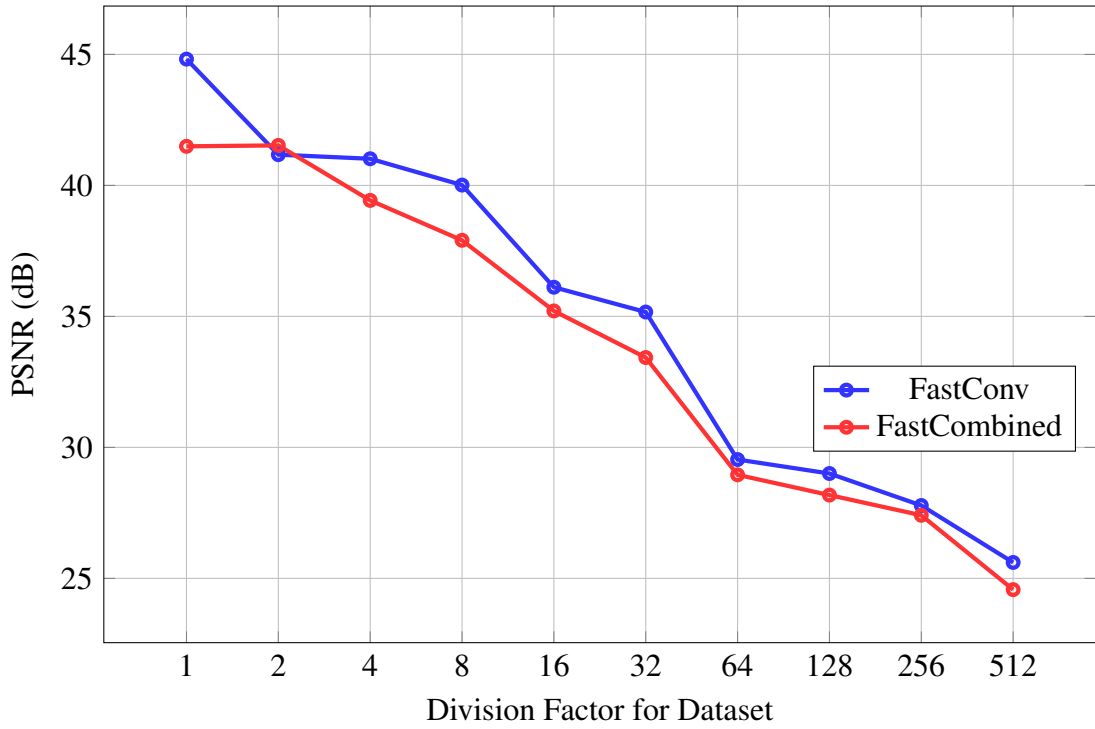


Figure 5.1.: This graph shows the PSNR for the model FastConv and FastCombined (see Table 5.1) for different subsections of the dataset. Only training and validation set sizes were reduced, all models were evaluated on the full test set.

factors, while between division factors 8 and 64 there is a high gradient in PSNR and spectral angle. A division factor of 8 on the HySpecNet-11k dataset results in 1004 training samples. Since the HySpecNet-11k dataset uses 202 spectral channels, we therefore require $1004/202 \approx 4.97$ training samples per channel until performance starts to decrease more quickly. Since for machine learning typically 5 training samples per dimension are recommended [57], we conclude that our models are not especially strongly effected by the Hughes Phenomenon. Additionally, the Combined Model approach does not increase the effect of the phenomenon. However, while the improvement of the model achieved by increasing the dataset size is reduced for dataset sizes above this threshold, the metrics still achieve the best values when training on the complete dataset. For this reason we continued to use the full dataset for the following experiments.

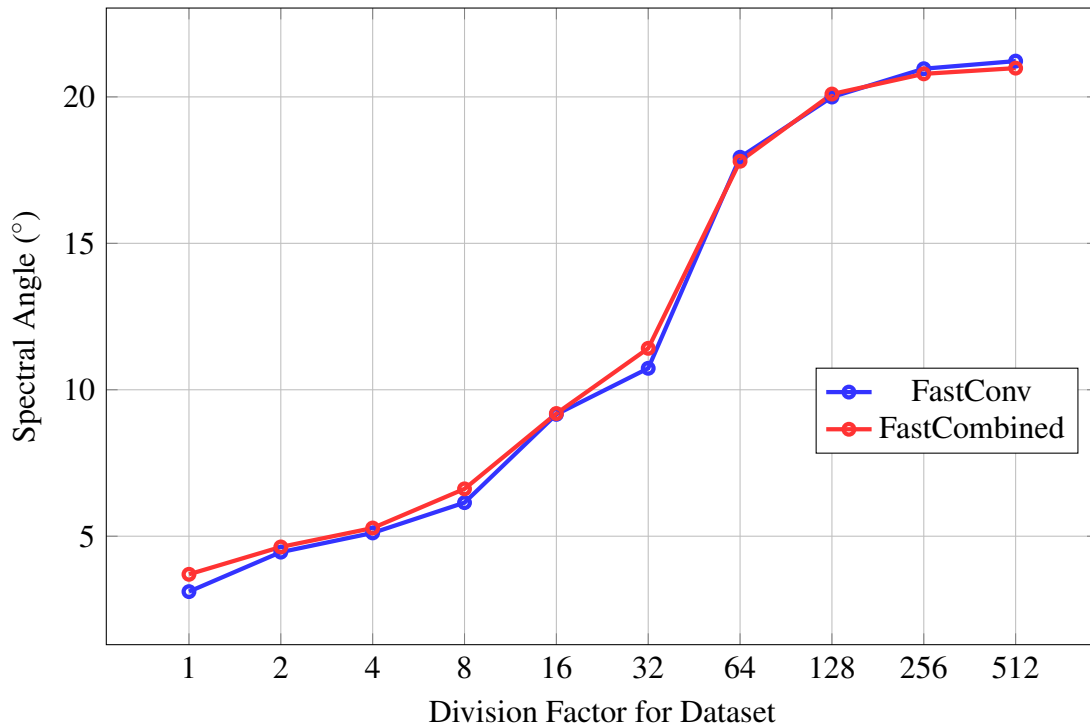


Figure 5.2.: This graph shows the spectral angle for the model FastConv and FastCombined (see Table 5.1) for different subsections of the dataset. Only training and validation set sizes were reduced, all models were evaluated on the full test set.

5.4. Analysis of Dual MSE Loss And Latent Spaces

For this experiment shown in Figure 5.3 we tested multiple different values for λ in the DualMSE loss. The spectral angle values for this experiment are shown in Figure A.1 in the appendix, as it shows similar results to Figure 5.3. Multiple variations of the Combined Model were tested in this way. All models in this experiment were trained by first pretraining the spectral autoencoder for 70 epochs and then training the complete Combined Model for 50 epochs. Furthermore all models had a spectral downsampling factor of 16 and a spatial downsampling factor of 4, resulting in a total compression ratio of 64 or a bitrate of 0.5 bpppc.

The tested models differ in two ways. Firstly, some models are of the **Combined** type while some are of the **FastCombined** type (see Table 5.1). The other variable in which the models differ is whether the weights of the spectral autoencoder are frozen after pretraining. We tested freezing all weights after pretraining, making the complete

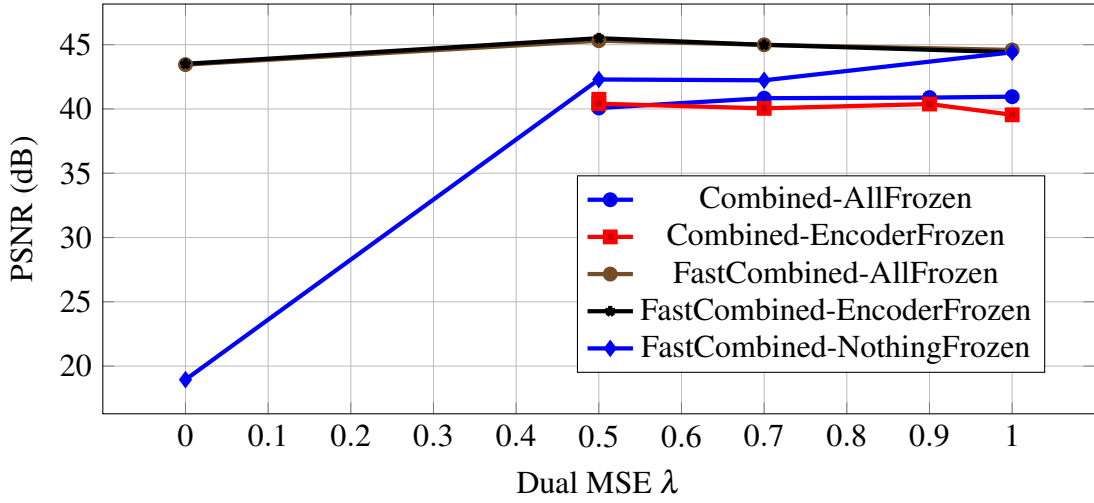


Figure 5.3.: This graph shows the PSNR for different values of the parameter λ for the DualMSE loss function.

model equivalent to training the spatial autoencoder while using the pretrained spectral encoder/decoder pair for pre- and postprocessing respectively. We also tested only freezing the weights of the spectral encoder after pretraining. In this way, the latent space produced by the spectral encoder does not change during training which we hoped could improve the training of the spatial autoencoder. Additionally, in contrast to the models where all weights of the spectral autoencoder are frozen after pretraining, the spectral decoder can still learn to adapt to the distortion introduced by the spatial autoencoder.

However, with the exception of $\lambda = 0$ for the model where only the encoder is frozen after pretraining, no clear trends in the PSNR or spectral angle could be observed regarding the value of λ . However, it is valuable to freeze either the spectral encoder or the complete spectral autoencoder, since the PSNR of the model **FastCombined-NothingFrozen** is 1.5 – 2 dB below the other **FastCombined** models for $\lambda = 0.5$ and $\lambda = 0.7$. For $\lambda = 0$ **FastCombined-NothingFrozen** is much worse with a difference of over 20 dB. To analyse why the model **FastCombined-NothingFrozen** with $\lambda = 0$ performed much worse than all other models, we computed the minimum and maximum values for the latent vectors y of the spectral encoder and for the reconstructed latents \hat{y} produced by the spatial autoencoder. For the model **FastCombined-EncoderFrozen** and **FastCombined-AllFrozen** these values behaved as expected with the maxima close to 1 and the minima close to 0. This is expected because the last layer of the spectral encoder is a sigmoid layer which constrains the latents to be in the interval $[0, 1]$. However, for the model **FastCombined-NothingFrozen** both the minimum and the maximum approached 0.5. On the test set the minimum latent value was 0.4875 and the maximum 0.5118. The

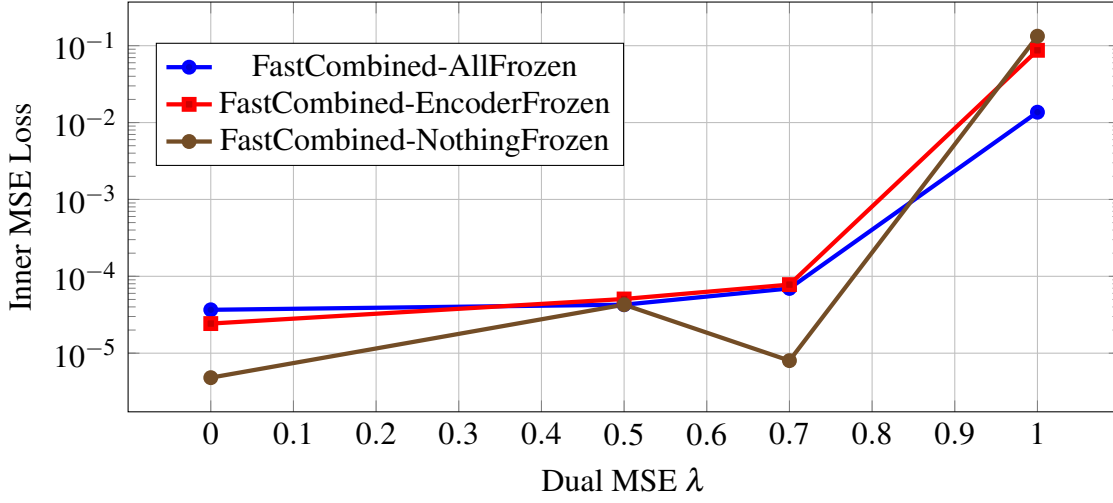


Figure 5.4.: This graph shows the Inner MSE Loss for different values of the parameter λ for the DualMSE loss function. Inner MSE Loss is the MSE between the output of the spectral encoder and the reconstruction of that latent from the spatial autoencoder in the Combined Model.

spatial autoencoder then reconstructed a constant output close to 0.5 for each channel. Because for $\lambda = 0$ only the MSE between y and the \hat{y} is relevant, the model managed to minimize this loss by learning to reduce the distance of the inputs between the latents.

However, while the usage of the DualMSE loss does not have a strong effect on PSNR or spectral angle with our models on the HySpecNet-11k dataset, it does succeed in strongly reducing the differences between the latent produced by the spectral encoder and its reconstruction. As shown in Figure 5.4, using DualMSE loss with $\lambda < 1$ drastically reduces the inner MSE, which is the MSE between y and \hat{y} . For **FastCombined-NothingFrozen** the inner MSE reduces from 0.133 to $8.00 \cdot 10^{-6}$, which is a factor of ≈ 16606 . For **FastCombined-EncoderFrozen** this factor is ≈ 1116 and for **FastCombined-AllFrozen** the reduction factor is ≈ 196 , which are still large decreases in loss. This can also be seen visually when comparing the reconstructed latents as we have done in Figure 5.5. We show the effect for **FastCombined-AllFrozen** since the frozen encoder gives the advantage of producing the same latent representations for all compared models, as the weights of the encoder are constant during training. It can be seen that with $\lambda = 1$ the latents have a visible grid pattern that is not present when using $\lambda < 1$. A similar pattern can also be seen during the first epoch of training for models with $\lambda < 1$ and is therefore likely a result of the convolutional layers used in the 2D CNN-based spatial autoencoder. During training these patterns disappear when the MSE between y and \hat{y} is part of the loss function, which is not the case for $\lambda = 1$. However, the models with $\lambda = 1$ produce

similar overall reconstructions \hat{x} as shown by the similar PSNR and spectral angle. The spectral decoder therefore is able to ignore these grid patterns even without additional training, since even the models **FastCombined-AllFrozen** and **Combined-AllFrozen** with $\lambda = 1$ have a high PSNR and a low spectral angle. This shows a high robustness of the spectral autoencoders with regard to distortions introduced to the bottleneck. It also demonstrates the high degree of explainability provided by the Combined Model architecture. The latent representations seen in Figure 5.5 also confirm the thesis that the 2D spectral encoder does preserve spatial dependencies. It can also be seen by the inner latents z in Figure 5.6 that the 2D CNN-based spatial autoencoder partially exploits the spatial dependencies in the latents as the inner latents have less spatial structure than the outer latents y . Especially for the second channel shown the much of the spatial structure seems to be exploited, while for the first channel the effect is not as strong. However, it can also be seen that some spatial dependencies still remain in z . For this reason an improved spatial autoencoder that is able to fully exploit the spatial dependencies in y could potentially improve the model.

The best performing model for this experiment was **FastCombined-EncoderFrozen** with $\lambda = 0.5$ with a PSNR of 45.50 dB and a spectral angle of 3.06° . The second best model was **FastCombined-AllFrozen** with $\lambda = 0.7$ with a PSNR of 45.31 and a spectral angle of 3.22° .

5.5. Analysis of Models Using Hyperprior-based Spatial Compression

In this experiment we compared multiple models using the Combined Model architecture with a hyperprior-based spatial autoencoder with each other as well as with the hyperprior-based spatial autoencoders outside of the Combined Model. For a detailed explanation of the hyperprior architecture see Section 3.4 and for its usage in the Combined Model see Subsection 4.3.3. The models **CombinedWithAttention** and **CombinedWithHyperprior** (see Table 5.1) are part of the first category, while Attention-Based and Scale Hyperprior are the respective spatial autoencoders tested separately from the Combined Model architecture without a reduction in the number of channels beforehand. As usual, we trained the Combined Model approaches by pretraining the spectral autoencoder for 70 epochs. We then froze both the spectral encoder and decoder while training the complete Combined Model for 50 additional epochs.

The baseline spatial hyperprior-based autoencoder models Attention-Based [42] and

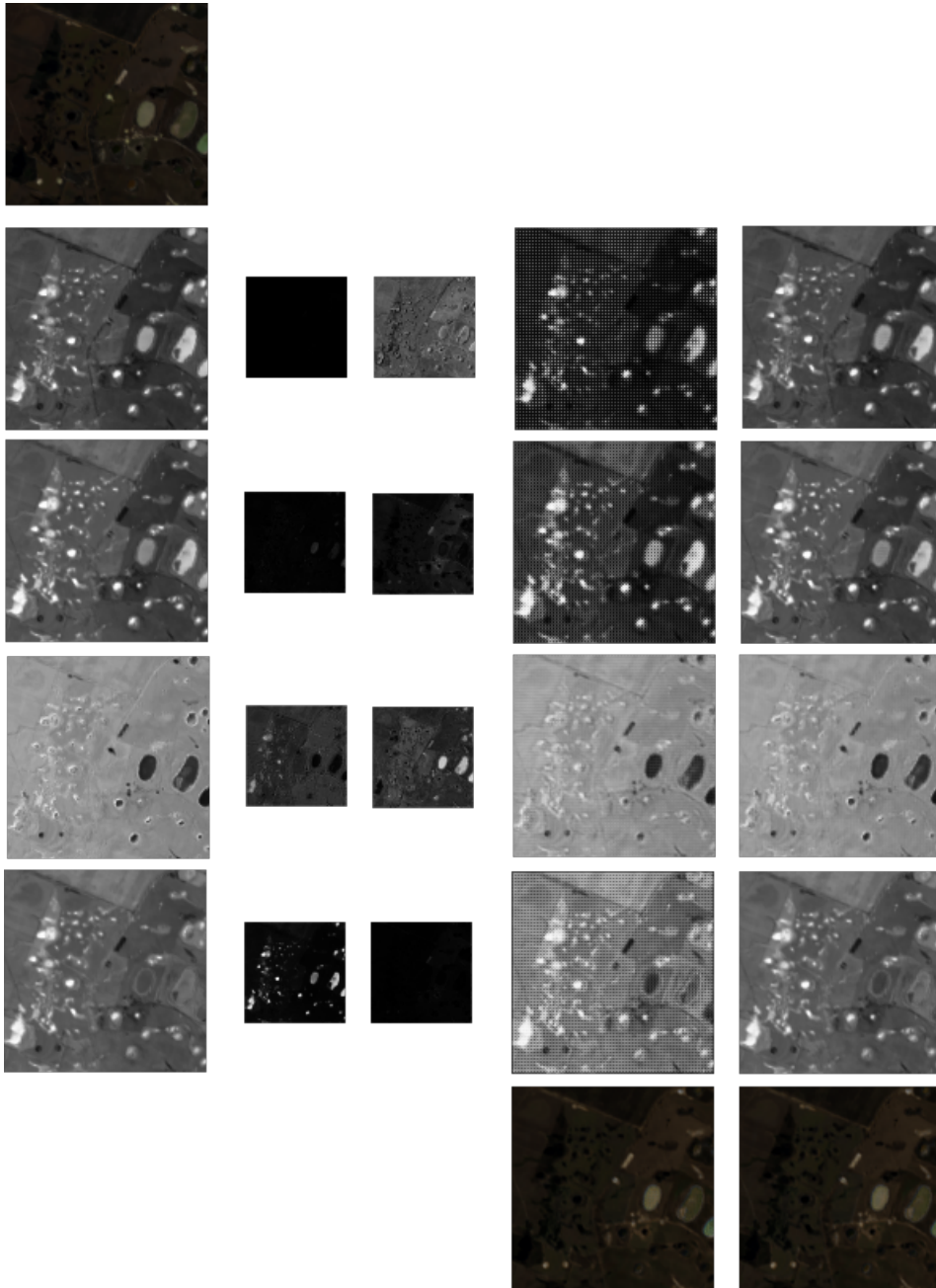


Figure 5.5.: This image shows four of the 13 channels in the latent representation of an image from the HySpecNet-11k dataset for the model FastCombined-AllFrozen for the values 1 and 0.7 for the parameter λ of the DualMSE loss function. The input is on the left. Below it is the latent produced by the spectral encoder, which is shared by both models. Following are, in order, the inner latent produced by the spectral encoder for $\lambda = 1$, and for $\lambda = 0.7$, the reconstruction of the latent produced by the spatial decoder for $\lambda = 1$ and for $\lambda = 0.7$. The full image reconstructions are shown on the bottom.

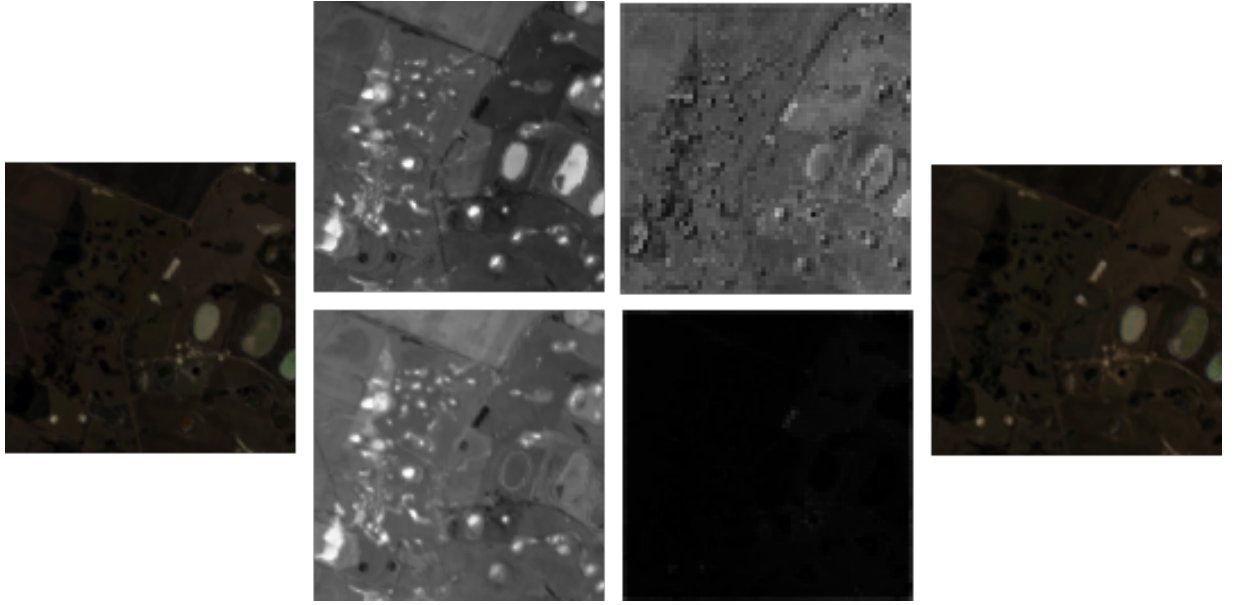


Figure 5.6.: Inner latent comparison

Table 5.2.: Comparison of hyperprior-based models. F is the spatial compression/decompression factor in the main encoder/decoder pair.

Model	F	Loss	φ	bpppc	PSNR (dB)	Spectral Angle ($^{\circ}$)
CombinedWithAttention	1	RD	0.1	0.0490	21.54	21.295
CombinedWithAttention	4	RD	0.1	0.0012	29.67	9.547
CombinedWithAttention	16	RD	0.1	0.0021	37.36	7.129
CombinedWithAttention	64	RD	0.1	0.0021	38.99	5.043
CombinedWithAttention	256	RD	0.1	0.0019	36.22	6.481
CombinedWithAttention	1	RD	5	8.3317	40.44	4.329
CombinedWithAttention	4	RD	5	0.0199	39.32	5.159
CombinedWithAttention	16	RD	5	0.0176	43.65	3.541
CombinedWithAttention	64	RD	5	0.0151	43.39	3.697
CombinedWithAttention	256	RD	5	0.0086	41.39	4.005
CombinedWithHyperprior	256	RD	0.1	0.0082	19.30	24.56
CombinedWithHyperprior	256	MSE	-	1.071	36.80	5.79
CombinedWithAttUnadjusted	16	RD	0.1	0.4327	36.58	7.503
CombinedWithAttUnadjusted	64	RD	0.1	0.4198	38.819	5.109
CombinedWithoutQuant	4	MSE	-	7.6030	48.01	2.816
CombinedWithAttention	4	MSE	-	2.9678	42.43	3.871
Attention-Based ([42])	256	RD	0.1	0.0012	33.92	8.71
Scale Hyperprior ([21])	256	RD	0.1	0.010	20.92	24.765

Scale Hyperprior [21] compress spatially by a factor of 16 in each spatial dimension in the main encoder. This leads to a total compression factor of $F=256$ in the main encoder. In order to compare the Combined Model approaches with the baselines, we also used a spatial compression factor of 256 in the main encoder. We compared these models using $\varphi = 0.1$. Our experiments showed that while the model **CombinedWithHyperprior** did not improve upon the baseline, the model **CombinedWithAttention** improved upon its baseline by 2.3 dB in PSNR and 2.22° in spectral angle. While it used a slightly higher bitrate of 0.0019 bpppc compared to the bitrate of the model Attention-Based of 0.0012 bpppc, both models have very low bitrates. For this reason, we focused on decreasing the distortion to achieve closer values to the non-hyperprior-based models while keeping a low bitrate. Furthermore, the attention-based hyperprior performed significantly better than the simple hyperprior-based model. This shows that the inclusion of self-attention modules as well as the context model significantly improves compression performance for a hyperprior-based hyperspectral image compression model. For this reason we used the model **CombinedWithAttention** for the following experiments.

We first modified the model by changing the spatial compression factor in the main encoder. Since the arithmetic coder compressed data losslessly, the quantized output of the main encoder is the direct input to the main decoder, which produces the reconstruction of the input. If we can reduce the loss of information in the spatial encoder, this could therefore reduce the amount of distortion present in the reconstruction. However, simply reducing the compression factor in the main part of the model without modifying the hyperprior autoencoder would increase the size of the bottleneck in the hyperprior autoencoder, therefore increasing the amount of side information that needs to be transmitted. Since the RD loss function tries to minimize bitrate, this can also lead to an increase in distortion. We therefore adjusted the spatial compression factor in the hyperprior autoencoder such that the size of the bottleneck in the hyperprior autoencoder is constant. We verified that this approach leads to both a lower distortion and a significantly lower bitrate for the spatial compression factors $F=64$ and $F=16$ with $\varphi = 0.1$. The model without adjusting the hyperprior autoencoder is the model **CombinedWithAttUnadjusted**. For $F=16$ the adjusted model achieves 2.41 dB PSNR more with a bitrate that is lower by a factor of $\frac{0.4327}{0.0021} \approx 206$. Additionally, we tested the models with $\varphi = 5$ and a selection of models using MSE loss. As explained in Subsection 5.2.2, MSE loss corresponds to the limit of MSE loss for $\varphi \rightarrow \infty$. We can therefore use MSE loss to find a lower bound for the possible amount of distortion achievable with a specific model using RD loss. In this way, we show that the model **CombinedWithHyperprior** performs worse than **CombinedWithAttention** independent of the value of φ by showing that **CombinedWithHyperprior** trained using MSE loss achieves 36.80 dB PSNR and a spectral angle of 5.79° , while **CombinedWithAttention** using the same spatial compression factor achieves a PSNR of 41.39 dB and a spectral angle of 4.005° for $\varphi = 5$.

Table 5.3.: Comparison of different kernel sizes for the spatial autoencoder and per-channel spatial autoencoder with the baseline model. For all models the spectral encoder is frozen after pretraining and $\lambda = 0.5$.

Kernel Size	Per Channel?	PSNR (dB)	Spectral Angle ($^{\circ}$)
3	No	45.50	3.06
5	No	43.72	3.30
7	No	42.59	3.58
3	Yes	42.33	3.30

The results of the experiments showed that the model **CombinedWithAttention** achieves the lowest distortion for spatial compression factors between $F=16$ and $F=64$ on the HySpecNet-11k dataset. For $\varphi = 0.1$ the lowest distortion of 38.99 dB PSNR and a spectral angle of 5.043° is achieved with a bitrate of 0.0021 bpppc for $F=64$. For $\varphi = 5$ the lowest distortion is 43.65 dB PSNR and a spectral angle of 3.541° and is achieved with $F=16$ with a bitrate of 0.0176 bpppc. While $F=256$ achieves slightly lower bitrates, this also increases distortion significantly. $F = 4$ results in significantly higher distortion without a significant decrease in bitrate, while $F=1$ both increases bitrate and distortion. With regard to the value of φ , a higher value of φ increases bitrate and distortion, as expected. The optimal value of φ therefore depends on the desired target bitrate or distortion.

Additionally, we performed an experiment to test the effect of the quantization step on distortion in the reconstruction in hyperprior-based compression models. **CombinedWithoutQuant** uses the same main autoencoder as the model **CombinedWithAttention**. However, it does not use an arithmetic coder or a hyperprior autoencoder in general. Additionally, it does not perform the quantization necessary for a hyperprior model. We compare this model with **CombinedWithAttention** also using $F=4$, trained with MSE loss. Since the model does not optimize for bitrate in this case and the arithmetic coder is lossless, the only difference between these models pertaining to distortion is the quantization step. **CombinedWithoutQuant** achieves a PSNR of 48.01 dB which is an improvement of 5.58 dB over the other model. This shows that the quantization step does increase distortion of hyperprior-based models significantly.

5.6. Analysis of Kernel Size Change in 2D CNN-based Spatial Autoencoder

To improve the amount of spatial dependencies exploited by the 2D CNN-based spatial autoencoder as part of the Combined Model, we experimented with changing the kernel sizes of the convolutional layers of the spatial autoencoder from three to the higher values five and seven. These experiments were performed using the model **FastCombined-EncoderFrozen** with $\lambda = 0.5$. However, as seen in Table 5.3, this did not improve either PSNR or spectral angle. The highest metrics were achieved by the base model using a kernel size of three with 45.50 dB and a spectral angle of 3.06° . The model using a kernel size of five achieved a lower PSNR of 43.72 dB and a higher spectral angle of 3.30° , the model using kernel size of seven resulted in the lowest acpsnr of 42.59 dB and the lowest spectral angle with 3.58° . We therefore conclude that a low kernel size of three is optimal for the CNN-based spatial autoencoder.

5.7. Analysis of Per-Channel 2D CNN-based Spatial Autoencoder

On the HySpecNet-11k dataset, the JPEG 2000 baseline performed surprisingly well. JPEG 2000 is executed per channel, therefore performing only spatial compression on 202 single-channel images. However, the 2D CNN-based spatial autoencoder uses the same filters to compress all channels in the bottleneck of the spectral autoencoder. We theorized that the differences between the latent channels might be too large to be appropriately captured using the same filters.

To test this thesis, we modified the 2D CNN-based spatial autoencoder to use a separate set of filters for each channel in the latent. Since this multiplies the number of parameters by the number of latent channels, we reduced the number of filters per channel. For the layers which originally had 256 filters, the modified model only uses 32 filters. For the layers that originally had 512 filters, the modified model uses 64 filters. The reduction in the number of filters is not only appropriate to reduce the number of parameters for the model, it follows also from the fact that less parameters are required for a convolutional layer that captures dependencies for a single channel than a convolutional layer that captures dependencies for multiple channel simultaneously. In total, the number of parameters in the spatial encoder after the reduction in the number of filters rises by 49.3 % from 50,849,818 to 75,928,346 parameters. We tested the modified model using the

model **FastCombined-EncoderFrozen** with $\lambda = 0.5$. However, the modification did not significantly change the performance of the Combined Model with regard to the metrics PSNR or the spectral angle as seen in Table 5.3.

5.8. Rate Distortion Comparison

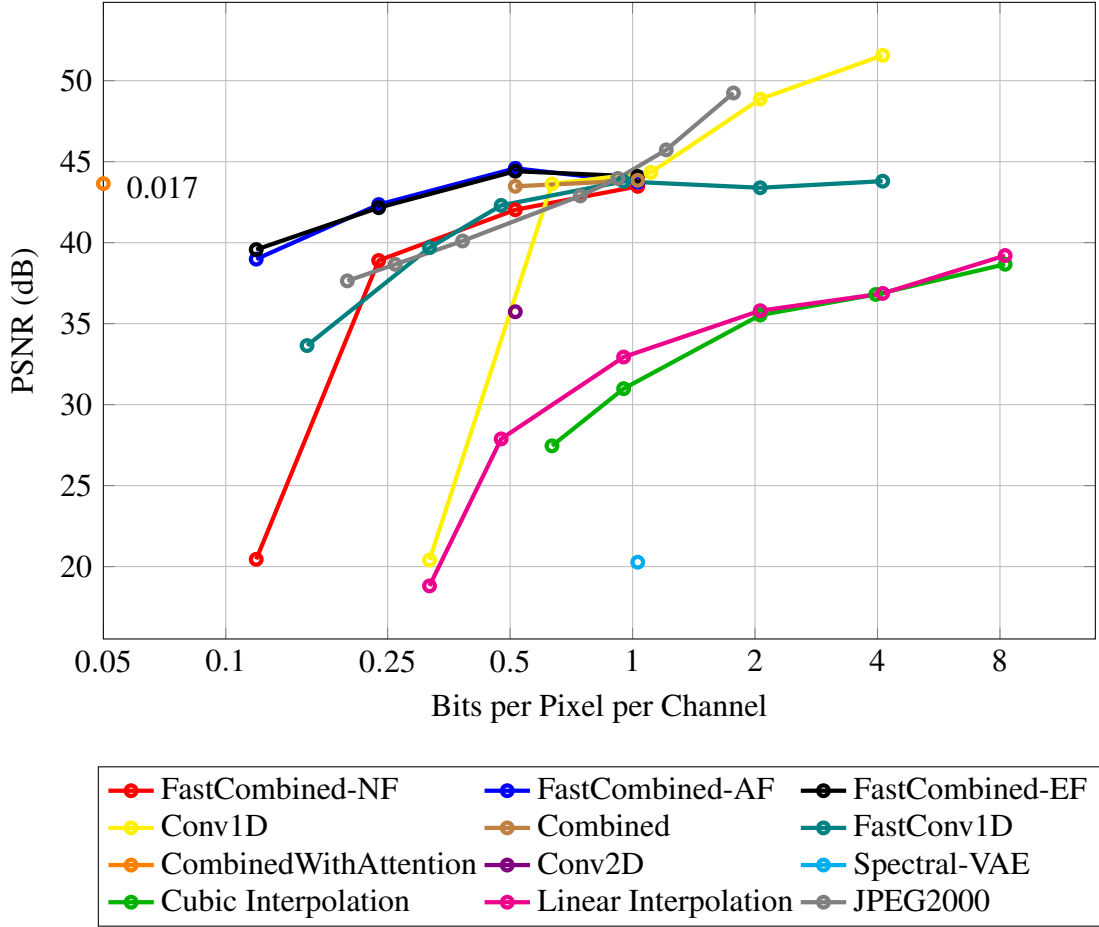


Figure 5.7.: Rate distortion plot showing PSNR for different bitrates for multiple model types.

Since the distortion produced by a lossy compression method depends on the compression rate [35], we compared the different model architectures used in this thesis for multiple bitrates. We compare all model types listed in Table 5.1 as well as multiple baselines. The models **Conv1D**, **FastConv1D** and **Spectral-VAE** were trained for 70 epochs. The models **Combined**, **FastCombined** and **CombinedWithAttention** were trained

for 50 epochs using the corresponding pre-trained spectral autoencoder. The model **Conv2D** was also trained for 50 epochs. For the model **FastCombined**, three variations were compared. For **FastCombined-NF**, the model is trained without freezing either the spectral encoder or decoder after pretraining, for **FastCombined-EF**, the spectral encoder is frozen after pretraining and for **FastCombined-AF**, the spectral encoder and decoder are frozen after pretraining. While the differences between these options were already discussed for a fixed bitrate of 0.5 bpppc in Section 5.4, they are also included in this experiment because of the surprising result that freezing the spectral encoder becomes significantly more important for lower bitrates, as discussed below in this section. The model **FastCombined-NF** uses DualMSE loss with $\lambda = 0.5$, the models **FastCombined-EF** and **FastCombined-AF** use $\lambda = 1$. For the models **Combined** and **CombinedWithAttention**, the best-performing options were used. For **Combined**, this is **EncoderFrozen** with $\lambda = 0.5$ and for **CombinedWithAttention** this is **AllFrozen** with $\lambda = 1$. The models **Combined** and **FastCombined** use a spatial compression factor of 4. The highest tested bitrate for such a model is therefore $\frac{1}{4}$ of the highest tested bitrate for the models **Conv1D** and **FastConv1D**.

Multiple baselines were used to quantify the performance of the models. The first two baselines, **Linear Interpolation** and **Cubic Interpolation** perform simple spectral interpolation. For compression using these methods, a given number of equidistant spectral channels are retained and all others discarded. For decompression the original channels are then restored using either linear or cubic interpolation. The first and last channel are always retained, since otherwise restoration of all channels would be impossible without extrapolating values. These methods were used as baselines because of the assumption that nearby spectral channels have similar values for the same pixel, since they capture wavelengths that are close together. The third baseline used was JPEG 2000, which was applied per-channel. The tested quality settings for JPEG 2000 were 1, 2, 3, 4, 5, 10, 15 and 20. We also consider **Conv1D** a baseline, since it is based on the model proposed in Kuester et.al. [23, 24].

The results show that the model **Spectral-VAE** is not successful in performing hyper-spectral image compression by reducing the number of spectral channels. This is not surprising since the VAE architecture in RGB image compression is also not commonly applied directly but rather through using a hyperprior architecture. A hyperprior-based spectral autoencoder is tested and discussed in Section 5.5. All other tested model outperform the interpolation-based baselines, showing that the spectral dependencies are too complex to be captured using only interpolation.

In comparison with the JPEG 2000 baseline and the models using the Combined Model architecture, the purely spectral compression methods **Conv1D** and **FastConv1D** perform

best for bitrates close to 0.5 bpppc. The model **Conv1D** outperforms JPEG 2000 for the bitrate 0.63 bpppc. Note that different model types can be tested for slightly different bitrates because of different padding strategies. Additionally, JPEG 2000 can only be set to a fixed number of different quality settings and therefore does not allow for setting a specific bitrate. The model **FastConv1D** has a higher PSNR than JPEG 2000 for bitrates from 0.31 to 0.47 bpppc. These correspond to two and respectively three channels in the latent representation. For higher bitrates, JPEG 2000 outperforms the spectral encoding methods. Further, the model **Conv1D** has a lower limit on the number of spectral channels in the latent dimensions, below which distortion strongly increases. This point is at 2 latent channels, which corresponds to a bitrate of 0.31 bpppc, where it achieves a PSNR of 20.40 dB, while for 4 latent channels, the model achieves 43.63 dB for a bitrate of 0.63 bpppc.

The Combined Model architecture outperforms all baselines and other compared models for bitrates at or below 0.5 bpppc. The model **Combined** has a PSNR of 43.48 when using 13 latent channels which results in a bitrate of 0.51 bpppc. For 26 latent channels, the model does not significantly improve, achieving a PSNR of 43.85 dB with a bitrate of 1.03 bpppc, which is worse than **Conv1D**, **FastConv1D** and JPEG 2000 for similar bitrates. This is expected, since the idea of the Combined Model architecture is to first reduce the number of spectral channels in order to be able to use spatial compression methods that work best for images with a low number of spectral channels such as RGB or multispectral images. This is also illustrated by the model **Conv2D**, which shows the 2D CNN-based spatial autoencoder directly applied to the hyperspectral data. This model only achieves a PSNR of 35.734 dB for a bitrate of 0.51 bpppc, which shows that a reduction in the amount of channels can strongly improve the performance of spatial compression methods. The effect is further shown by the model **AttentionHyperprior**, which uses the spatial autoencoder of the **CombinedWithAttention** model without first reducing the number of spectral channels, achieving a much lower PSNR of 33.92 dB. The model **FastCombined-AF** achieves the highest PSNR of all compared models, 44.59 dB, for the bitrate 0.51. For the bitrates 0.23 bpppc and 0.12 bpppc it is the second best model behind **CombinedWithAttention** with PSNRs of 42.35 and 38.98 respectively. However, the model **FastCombined-NF** performs much worse for the bitrate 0.12 with only a PSNR of 20.44 dB. It follows that freezing the spectral encoder is important when training for very low bitrates.

For all bitrates below 0.25, the highest PSNR is achieved by the model **CombinedWithAttention**. Additionally, this distortion is achieved with a much lower bitrate than in other models with similar PSNR values. In comparison with the second model with the second lowest bitrate for similar distortion, **FastCombined-AF**, the compression rate differs by a factor of $0.24/0.017 = 14.11$. We conclude that for the HySpecNet-11k the

CombinedWithAttention model is optimal in order to achieve bitrates < 0.25 bpppc, while the **FastCombined** model can achieve slightly lower distortion for bitrates ≈ 0.5 bpppc. Additionally, the results show that the Combined Model approach improves upon the purely spectral compression methods and purely spatial compression for bitrates lower than 1 bpppc.

6. Discussion and Conclusion

This thesis investigates the capabilities of utilising a spatio-spectral approach to learning-based hyperspectral image compression. To this end, a novel architecture is proposed that combines a spectral and a spatial autoencoder by applying the spatial autoencoder to the compression bottleneck of the spectral autoencoder. We call this architecture Combined Model. Multiple different approaches are used for both the spectral and spatial parts of the Combined Model. Additionally, a method of training optimised for the Combined Model architecture is proposed. This method employs pretraining of the spectral autoencoder in conjunction with freezing parts of the spectral autoencoder during actual training as well as a novel loss function called DualMSE loss to improve the capabilities of the Combined Model.

For spectral encoding we propose two main model types. The first method is a 1D CNN model that performs compression in the spectral dimension separately per pixel of an input image. We modified the model to allow for variable compression ratios. Furthermore, we adapted the model to the dataset used in this thesis by changing the number of filters in each layer and changed the padding strategy to be more compatible with a variable compression ratio. The second method is a 2D CNN model that employs a kernel size of 1 in the convolutional layers to preserve spatial dependencies in the input while allowing for faster training and full control over the size of the latent dimension. Thirdly, a spectral variational autoencoder is also proposed. However, because experiments show that this model is not viable for spectral autoencoding. Therefore the focus of this thesis is on the other two spectral compression models.

For spatial encoding three main architectures are proposed. The first is a 2D CNN-based model that reduces the spatial dimensionality of the input by a fixed compression ratio. This model is designed to have low complexity to be a baseline for the other spatial compression architectures as well as allow for testing of the Combined Model architecture and the proposed training method. The second spatial encoding architecture is a gaussian scale hyperprior that serves as a second baseline to the third spatial encoding architecture, which also employs a hyperprior. The third spatial encoding architecture is a state-of-the-art hyperprior-based model using a gaussian mixture model hyperprior, a context model and self-attention modules. Additionally, the model is adapted to the

Combined Model architecture by reducing the amount of spatial dependencies exploited by the main autoencoder and consequently increasing the amount of spatial dependencies exploited in the hyperprior autoencoder.

The performance of the different model architectures, specific modifications to the models and the training method designed for the Combined Model was tested on the HySpecNet-11k dataset, which provides large-scale hyperspectral image data of the surface of the earth with high spectral resolution. The models were evaluated on the distortion metrics PSNR and spectral angle as well as the achieved compression bitrate.

We conducted experiments to assess the effect of the Hughes Phenomenon on the Combined Model which showed that the effect on both the Combined Model and the 2D CNN-based spectral autoencoder are in line with typical expectations for machine learning tasks. With regard to the specific training method, we showed that pretraining the spectral autoencoder and freezing the spectral encoder during training leads to significant improvements in PSNR and spectral angle for a fixed bitrate. We further showed that the effect of freezing the spectral encoder during training is increased when the latent dimension of the spectral encoder is small. In addition to this, while the use of the novel DualMSE loss does not significantly decrease distortion, it is successful in reducing the distortion between the latent of the spectral encoder and its reconstruction obtained after applying the spatial autoencoder to that latent by multiple orders of magnitude. The DualMSE loss is therefore potentially useful for an application of the Combined Model architecture that relies on a latent reconstruction with low distortion. To achieve this result, we exploited that the Combined Model architecture allows for a high degree of explainability by studying the spectral latents. Experiments also showed that the Combined Model achieves significantly lower distortion than all compared benchmarks on bitrates at or below 0.5 bpppc. We thereby showed that exploiting both spatial and spectral dependencies of hyperspectral images is essential to improve the performance of learning-based hyperspectral compression methods on the HySpecNet-11k dataset. Further, the Combined Model significantly improves the use of hyperprior-based models on hyperspectral data, achieving lower distortion than hyperprior-based models without the Combined Model approach and significantly lower bitrates at comparative distortion levels than all models not using a hyperprior. The Combined Model with the hyperprior-based spatial encoder using self-attention modules is to our knowledge the only learning-based hyperspectral image compression model able to achieve very low bitrates (< 0.02 bpppc) with low distortion. We therefore propose new state of the art learning-based hyperspectral compression models for low and for very low bitrates.

We see multiple possible areas of future research that could be explored to improve on the results from this thesis. The latents of the CNN-based spatial autoencoder show that there

are still spatial dependencies present that could potentially be exploited by a different spatial autoencoder architecture. Additionally, while different spectral compression ratios were extensively tested, more research could be performed into different spatial compression ratios. While the hyperprior-based spatial autoencoder achieves this through the use of the arithmetic coder, it requires a quantization step in order to apply the arithmetic coder, which leads to some amount of irreversible loss during the encoding process. We showed in the experiments removing the quantization step significantly reduce distortion. However, since the arithmetic autoencoder is then unable to be used, this is not directly for compression. A novel way to perform the quantization with less loss of information or an approach to possibly not require quantization at all could significantly improve the hyperprior-based models. Improvements in this area could also potentially translate to the application of hyperprior-based RGB image compression. Furthermore, experiments were only performed on the HySpecNet-11k dataset. Research into the application of the models proposed in this thesis on different domains could provide more insight into the generalisability of the models proposed in this thesis.

Bibliography

- [1] M. Zea, A. Souza, Y. Yang, L. Lee, K. Nemali, and L. Hoagland, “Leveraging high-throughput hyperspectral imaging technology to detect cadmium stress in two leafy green crops and accelerate soil remediation efforts”, *Environmental Pollution*, vol. 292, p. 118 405, Jan. 1, 2022, ISSN: 0269-7491. DOI: 10.1016/j.envpol.2021.118405.
- [2] M. L. Henriksen, C. B. Karlsen, P. Klarskov, and M. Hinge, “Plastic classification via in-line hyperspectral camera analysis and unsupervised machine learning”, *Vibrational Spectroscopy*, vol. 118, p. 103 329, Jan. 1, 2022, ISSN: 0924-2031. DOI: 10.1016/j.vibspec.2021.103329.
- [3] L. Guanter *et al.*, “The EnMAP spaceborne imaging spectroscopy mission for earth observation”, *Remote Sensing*, vol. 7, no. 7, pp. 8830–8857, Jul. 2015, Number: 7 Publisher: Multidisciplinary Digital Publishing Institute, ISSN: 2072-4292. DOI: 10.3390/rs70708830.
- [4] N. Bohn *et al.*, “Glacier ice surface properties in south-west greenland ice sheet: First estimates from PRISMA imaging spectroscopy data”, *Journal of Geophysical Research: Biogeosciences*, vol. 127, no. 3, e2021JG006718, 2022, eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1029/2021JG006718>, ISSN: 2169-8961. DOI: 10.1029/2021JG006718.
- [5] A. B. Pascual-Ventoe *et al.*, “Prototyping crop traits retrieval models for CHIME: Dimensionality reduction strategies applied to PRISMA data”, *Remote Sensing*, vol. 14, no. 10, p. 2448, Jan. 2022, Number: 10 Publisher: Multidisciplinary Digital Publishing Institute, ISSN: 2072-4292. DOI: 10.3390/rs14102448.
- [6] V. Döpper *et al.*, “Estimating soil moisture content under grassland with hyperspectral data using radiative transfer modelling and machine learning”, *International Journal of Applied Earth Observation and Geoinformation*, vol. 110, p. 102 817, Jun. 1, 2022, ISSN: 1569-8432. DOI: 10.1016/j.jag.2022.102817.
- [7] R. Loizzo *et al.*, “Prisma mission status and perspective”, in *IGARSS 2019 - 2019 IEEE International Geoscience and Remote Sensing Symposium*, ISSN: 2153-7003, Jul. 2019, pp. 4503–4506. DOI: 10.1109/IGARSS.2019.8899272.

-
- [8] R. Guarini, R. Loizzo, F. Longo, S. Mari, T. Scopa, and G. Varacalli, "Overview of the prisma space and ground segment and its hyperspectral products", in *2017 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, ISSN: 2153-7003, Jul. 2017, pp. 431–434. DOI: 10.1109/IGARSS.2017.8126986.
- [9] R. Guarini *et al.*, "Prisma hyperspectral mission products", in *IGARSS 2018 - 2018 IEEE International Geoscience and Remote Sensing Symposium*, ISSN: 2153-7003, Jul. 2018, pp. 179–182. DOI: 10.1109/IGARSS.2018.8517785.
- [10] N. Acito, M. Diani, and G. Corsini, "PRISMA spatial resolution enhancement by fusion with sentinel-2 data", *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 15, pp. 62–79, 2022, Conference Name: IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, ISSN: 2151-1535. DOI: 10.1109/JSTARS.2021.3132135.
- [11] M. Drusch *et al.*, "Sentinel-2: ESA's optical high-resolution mission for GMES operational services", *Remote Sensing of Environment*, The Sentinel Missions - New Opportunities for Science, vol. 120, pp. 25–36, May 15, 2012, ISSN: 0034-4257. DOI: 10.1016/j.rse.2011.11.026.
- [12] D. Spiller, S. Amici, and L. Ansalone, "Transfer learning analysis for wildfire segmentation using prisma hyperspectral imagery and convolutional neural networks", in *2022 12th Workshop on Hyperspectral Imaging and Signal Processing: Evolution in Remote Sensing (WHISPERS)*, ISSN: 2158-6276, Sep. 2022, pp. 1–5. DOI: 10.1109/WHISPERS56178.2022.9955054.
- [13] D. Spiller, K. Thangavel, S. T. Sasidharan, S. Amici, L. Ansalone, and R. Sabatini, "Wildfire segmentation analysis from edge computing for on-board real-time alerts using hyperspectral imagery", in *2022 IEEE International Conference on Metrology for Extended Reality, Artificial Intelligence and Neural Engineering (MetroXRaine)*, Oct. 2022, pp. 725–730. DOI: 10.1109/MetroXRaine54828.2022.9967553.
- [14] C. E. Shannon, "A mathematical theory of communication", *Bell System Technical Journal*, vol. 27, no. 3, pp. 379–423, 1948, ISSN: 1538-7305. DOI: 10.1002/j.1538-7305.1948.tb01338.x.
- [15] D. J. C. MacKay, *Information theory, inference, and learning algorithms*. Cambridge, UK ; New York: Cambridge University Press, 2003, 628 pp., ISBN: 978-0-521-64298-9.
- [16] F. Garcia-Vilchez *et al.*, "On the impact of lossy compression on hyperspectral image classification and unmixing", *IEEE Geoscience and Remote Sensing Letters*, vol. 8, no. 2, pp. 253–257, Mar. 2011, ISSN: 1545-598X, 1558-0571. DOI: 10.1109/LGRS.2010.2062484.

- [17] E. Delp and O. Mitchell, “Image compression using block truncation coding”, *IEEE Transactions on Communications*, vol. 27, no. 9, pp. 1335–1342, Sep. 1979, Conference Name: IEEE Transactions on Communications, ISSN: 1558-0857. DOI: 10.1109/TCOM.1979.1094560.
- [18] M. M. Hannuksela, J. Lainema, and V. K. Malamil Vadakital, “The high efficiency image file format standard [standards in a nutshell]”, *IEEE Signal Processing Magazine*, vol. 32, no. 4, pp. 150–156, Jul. 2015, Conference Name: IEEE Signal Processing Magazine, ISSN: 1558-0792. DOI: 10.1109/MSP.2015.2419292.
- [19] S. Ruder, *An overview of gradient descent optimization algorithms*, Jun. 15, 2017. DOI: 10.48550/arXiv.1609.04747. arXiv: 1609.04747 [cs].
- [20] J. Ballé, V. Laparra, and E. P. Simoncelli, *End-to-end optimized image compression*, Number: arXiv:1611.01704, Mar. 3, 2017. DOI: 10.48550/arXiv.1611.01704. arXiv: 1611.01704 [cs, math].
- [21] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston, *Variational image compression with a scale hyperprior*, Number: arXiv:1802.01436, May 1, 2018. DOI: 10.48550/arXiv.1802.01436. arXiv: 1802.01436 [cs, eess, math].
- [22] D. Minnen, J. Ballé, and G. D. Toderici, “Joint autoregressive and hierarchical priors for learned image compression”, in *Advances in Neural Information Processing Systems*, vol. 31, Curran Associates, Inc., 2018.
- [23] J. 1. Kuester, W. 1. Gross, W. 1. 1. F. I. Middelmann, G. Image Exploitation, E. Fraunhofer IOSB, and G. Image Exploitation, “1d-convolutional autoencoder based hyperspectral data compression”, ISSN: 16821750 Num Pages: 15-21, Copernicus GmbH, 2021, pp. 15–21. DOI: 10.5194/isprs-archives-XLIII-B1-2021-15-2021.
- [24] J. Kuester, W. Gross, S. Schreiner, M. Heizmann, and W. Middelmann, “Transferability of convolutional autoencoder model for lossy compression to unknown hyperspectral prisma data”, in *2022 12th Workshop on Hyperspectral Imaging and Signal Processing: Evolution in Remote Sensing (WHISPERS)*, ISSN: 2158-6276, Sep. 2022, pp. 1–5. DOI: 10.1109/WHISPERS56178.2022.9955109.
- [25] R. La Grassa, C. Re, G. Cremonese, and I. Gallo, “Hyperspectral data compression using fully convolutional autoencoder”, *Remote Sensing*, vol. 14, no. 10, p. 2472, Jan. 2022, Number: 10 Publisher: Multidisciplinary Digital Publishing Institute, ISSN: 2072-4292. DOI: 10.3390/rs14102472.
- [26] Y. Guo, Y. Chong, Y. Ding, S. Pan, and X. Gu, “Learned hyperspectral compression using a student’s t hyperprior”, *Remote Sensing*, vol. 13, no. 21, p. 4390, Jan. 2021, Number: 21 Publisher: Multidisciplinary Digital Publishing Institute, ISSN: 2072-4292. DOI: 10.3390/rs13214390.

- [27] Y. Dua, V. Kumar, and R. S. Singh, “Comprehensive review of hyperspectral image compression algorithms”, *Optical Engineering*, vol. 59, no. 9, p. 090 902, Sep. 2020, Publisher: SPIE, ISSN: 0091-3286, 1560-2303. DOI: 10.1117/1.OE.59.9.090902.
- [28] A. Aidini, M. Giannopoulos, A. Pentari, K. Fotiadou, and P. Tsakalides, “Hyperspectral image compression and super-resolution using tensor decomposition learning”, in *2019 53rd Asilomar Conference on Signals, Systems, and Computers*, Pacific Grove, CA, USA: IEEE, Nov. 2019, pp. 1369–1373, ISBN: 978-1-72814-300-2. DOI: 10.1109/IEEECONF44664.2019.9048735.
- [29] C. Deng, Y. Cen, and L. Zhang, “Learning-based hyperspectral imagery compression through generative neural networks”, *Remote Sensing*, vol. 12, no. 21, p. 3657, Jan. 2020, Number: 21 Publisher: Multidisciplinary Digital Publishing Institute, ISSN: 2072-4292. DOI: 10.3390/rs12213657.
- [30] S. Kumar, S. Chaudhuri, B. Banerjee, and F. Ali, “Onboard hyperspectral image compression using compressed sensing and deep learning”, in *Computer Vision – ECCV 2018 Workshops*, L. Leal-Taixé and S. Roth, Eds., vol. 11130, Series Title: Lecture Notes in Computer Science, Cham: Springer International Publishing, 2019, pp. 30–42, ISBN: 978-3-030-11011-6. DOI: 10.1007/978-3-030-11012-3_3.
- [31] J. A. Saghri, “Adaptive two-stage karhunen-loeve-transform scheme for spectral decorrelation in hyperspectral bandwidth compression”, *Optical Engineering*, vol. 49, no. 5, p. 057 001, May 1, 2010, ISSN: 0091-3286. DOI: 10.1117/1.3425656.
- [32] A. Karami, M. Yazdi, and G. Mercier, “Compression of hyperspectral images using discrete wavelet transform and tucker decomposition”, *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 5, no. 2, pp. 444–450, Apr. 2012, Conference Name: IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, ISSN: 2151-1535. DOI: 10.1109/JSTARS.2012.2189200.
- [33] M. Conoscenti, R. Coppola, and E. Magli, “Constant SNR, rate control, and entropy coding for predictive lossy hyperspectral image compression”, *IEEE Transactions on Geoscience and Remote Sensing*, vol. 54, no. 12, pp. 7431–7441, Dec. 2016, Conference Name: IEEE Transactions on Geoscience and Remote Sensing, ISSN: 1558-0644. DOI: 10.1109/TGRS.2016.2603998.
- [34] M. Hernández-Cabronero *et al.*, “The CCSDS 123.0-b-2 “low-complexity lossless and near-lossless multispectral and hyperspectral image compression” standard: A comprehensive review”, *IEEE Geoscience and Remote Sensing Magazine*, vol. 9, no. 4, pp. 102–119, Dec. 2021, Conference Name: IEEE Geoscience and Remote Sensing Magazine, ISSN: 2168-6831. DOI: 10.1109/MGRS.2020.3048443.

- [35] T. Berger, “Rate-distortion theory”, in *Wiley Encyclopedia of Telecommunications*, eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/0471219282.eot142>, John Wiley & Sons, Ltd, 2003, ISBN: 978-0-471-21928-6. DOI: 10.1002/0471219282.eot142.
- [36] M. H. P. Fuchs and B. Demir, *HySpecNet-11k: A large-scale hyperspectral dataset for benchmarking learning-based hyperspectral image compression methods*, Jun. 2, 2023. DOI: 10.48550/arXiv.2306.00385. arXiv: 2306.00385 [cs, eess].
- [37] H. Shen, W. D. Pan, Y. Dong, and Z. Jiang, “Golomb-rice coding parameter learning using deep belief network for hyperspectral image compression”, in *2017 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, Fort Worth, TX: IEEE, Jul. 2017, pp. 2239–2242, ISBN: 978-1-5090-4951-6. DOI: 10.1109/IGARSS.2017.8127434.
- [38] D. Hong *et al.*, “SpectralFormer: Rethinking hyperspectral image classification with transformers”, *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–15, 2022, ISSN: 0196-2892, 1558-0644. DOI: 10.1109/TGRS.2021.3130716. arXiv: 2107.02988 [cs].
- [39] J. E. Fowler and J. T. Rucker, “Three-dimensional wavelet-based compression of hyperspectral imagery”, in *Hyperspectral Data Exploitation*, Section: 14 eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/9780470124628.ch14>, John Wiley & Sons, Ltd, 2007, pp. 379–407, ISBN: 978-0-470-12462-8. DOI: 10.1002/9780470124628.ch14.
- [40] Y. Hu, W. Yang, and J. Liu, “Coarse-to-fine hyper-prior modeling for learned image compression”, *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 7, pp. 11 013–11 020, Apr. 3, 2020, Number: 07, ISSN: 2374-3468. DOI: 10.1609/aaai.v34i07.6736.
- [41] M. J. Wainwright and E. P. Simoncelli, “Scale mixtures of gaussians and the statistics of natural images”, 1999.
- [42] Z. Cheng, H. Sun, M. Takeuchi, and J. Katto, *Learned image compression with discretized gaussian mixture likelihoods and attention modules*, Mar. 30, 2020. DOI: 10.48550/arXiv.2001.01568. arXiv: 2001.01568 [eess].
- [43] A. Vaswani *et al.*, *Attention is all you need*, Dec. 5, 2017. DOI: 10.48550/arXiv.1706.03762. arXiv: 1706.03762 [cs].
- [44] H. Liu *et al.*, *Non-local attention optimized deep image compression*, version: 1, Apr. 22, 2019. DOI: 10.48550/arXiv.1904.09757. arXiv: 1904.09757 [cs, eess].

- [45] M. Li, W. Zuo, S. Gu, D. Zhao, and D. Zhang, *Learning convolutional networks for content-weighted image compression*, Sep. 19, 2017. DOI: 10.48550/arXiv.1703.10553. arXiv: 1703.10553[cs].
- [46] F. Mentzer, E. Agustsson, M. Tschannen, R. Timofte, and L. Van Gool, *Conditional probability models for deep image compression*, Jun. 4, 2019. DOI: 10.48550/arXiv.1801.04260. arXiv: 1801.04260[cs].
- [47] K. O’Shea and R. Nash, “An introduction to convolutional neural networks”, *An Introduction to Convolutional Neural Networks*, Nov. 26, 2015, Publisher: arXiv.
- [48] J. Gu *et al.*, “Recent advances in convolutional neural networks”, *Pattern Recognition*, vol. 77, pp. 354–377, May 1, 2018, ISSN: 0031-3203. DOI: 10.1016/j.patcog.2017.10.013.
- [49] I. H. Witten, R. M. Neal, and J. G. Cleary, “Arithmetic coding for data compression”, *Communications of the ACM*, vol. 30, no. 6, pp. 520–540, Jun. 1987, ISSN: 0001-0782, 1557-7317. DOI: 10.1145/214762.214771.
- [50] A. Said, *Introduction to arithmetic coding – theory and practice*, Feb. 1, 2023. DOI: 10.48550/arXiv.2302.00819. arXiv: 2302.00819[cs,math].
- [51] V. Sze and M. Budagavi, “High throughput CABAC entropy coding in HEVC”, *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 22, no. 12, pp. 1778–1791, Dec. 2012, Conference Name: IEEE Transactions on Circuits and Systems for Video Technology, ISSN: 1558-2205. DOI: 10.1109/TCSVT.2012.2221526.
- [52] A. v. d. Oord, N. Kalchbrenner, O. Vinyals, L. Espeholt, A. Graves, and K. Kavukcuoglu, *Conditional image generation with PixelCNN decoders*, Jun. 18, 2016. DOI: 10.48550/arXiv.1606.05328. arXiv: 1606.05328[cs].
- [53] C. Szegedy *et al.*, *Going deeper with convolutions*, Sep. 16, 2014. DOI: 10.48550/arXiv.1409.4842. arXiv: 1409.4842[cs].
- [54] J. L. Ba, J. R. Kiros, and G. E. Hinton, *Layer normalization*, Jul. 21, 2016. DOI: 10.48550/arXiv.1607.06450. arXiv: 1607.06450[cs,stat].
- [55] F. A. Kruse *et al.*, “The spectral image processing system (SIPS) — interactive visualization and analysis of imaging spectrometer data”, *Remote Sensing of Environment*, Airbone Imaging Spectrometry, vol. 44, no. 2, pp. 145–163, May 1, 1993, ISSN: 0034-4257. DOI: 10.1016/0034-4257(93)90013-N.
- [56] G. Hughes, “On the mean accuracy of statistical pattern recognizers”, *IEEE Transactions on Information Theory*, vol. 14, no. 1, pp. 55–63, Jan. 1968, Conference Name: IEEE Transactions on Information Theory, ISSN: 1557-9654. DOI: 10.1109/TIT.1968.1054102.

-
- [57] S. Theodoridis and K. Koutroumbas, *Pattern recognition*, 4th ed. Burlington, MA London: Academic Press, 2009, ISBN: 978-1-59749-272-0.

A. Appendix

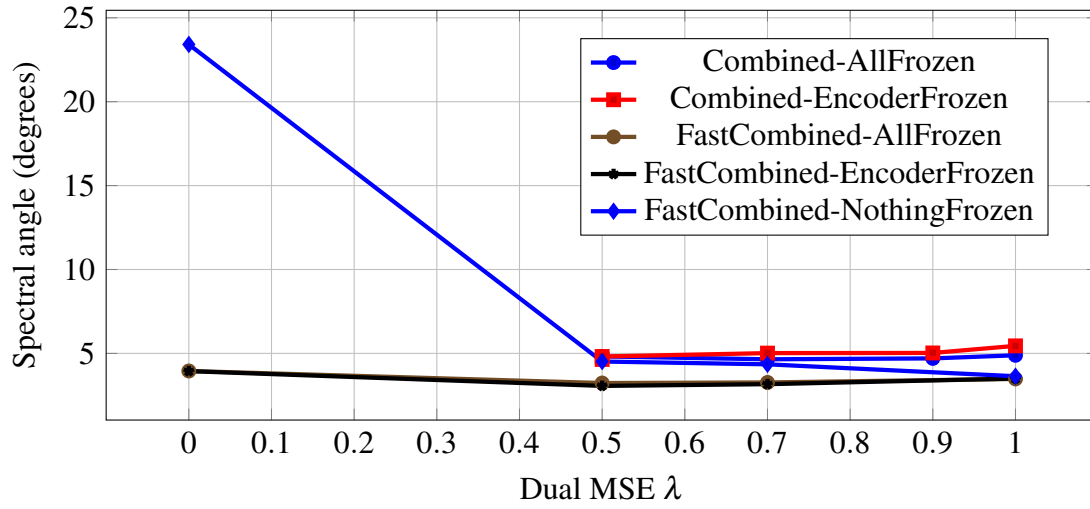


Figure A.1.: Dual MSE lambda plot

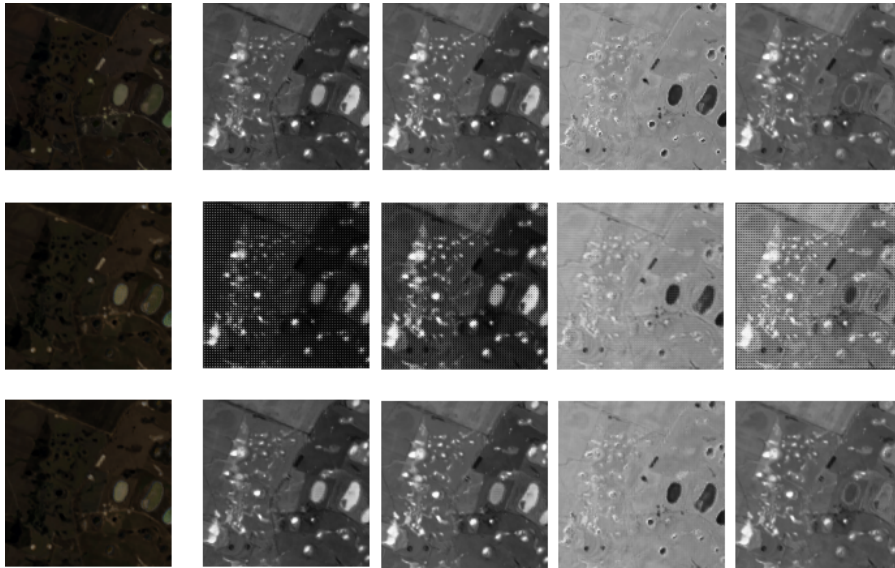


Figure A.2.: Latent image comparison