

Master Thesis



Czech
Technical
University
in Prague

F3

Faculty of Electrical Engineering
Department of Radioelectronics

Machine Learning Techniques for High Performance Image Compression

Metody strojového učení pro efektivní kompresi
obrazu

Rudolf Studený

Supervisor: Ing. Fliegel Karel Ph.D.
August 2021

I. Personal and study details

Student's name: **Studený Rudolf** Personal ID number: **434707**
Faculty / Institute: **Faculty of Electrical Engineering**
Department / Institute: **Department of Radioelectronics**
Study program: **Electronics and Communications**
Specialisation: **Audiovisual and Signal Processing**

II. Master's thesis details

Master's thesis title in English:

Machine Learning Techniques for High Performance Image Compression

Master's thesis title in Czech:

Metody strojového učení pro efektivní kompresi obrazu

Guidelines:

Give an overview of recent methods for image compression using machine learning techniques (ML). Focus on the latest applications of deep learning (DL) methods. For the selected problem and test data, analyze the performance of the machine learning based methods compared to classical approaches.

Bibliography / sources:

[1] Chen, Z., He, T., Learning based Facial Image Compression with semantic fidelity metric, Neurocomputing, 338, 2019.
[2] Minnen, D., Ballé, J., Toderici, G., Joint autoregressive and hierarchical priors for learned image compression, Advances in Neural Information Processing Systems, 2018.

Name and workplace of master's thesis supervisor:

Ing. Karel Fliegel, Ph.D., Department of Radioelectronics, FEE

Name and workplace of second master's thesis supervisor or consultant:

Date of master's thesis assignment: **14.02.2020** Deadline for master's thesis submission: **13.08.2021**

Assignment valid until: **30.09.2021**

Ing. Karel Fliegel, Ph.D.
Supervisor's signature

doc. Ing. Josef Dobeš, CSc.
Head of department's signature

prof. Mgr. Petr Páta, Ph.D.
Dean's signature

III. Assignment receipt

The student acknowledges that the master's thesis is an individual work. The student must produce his thesis without the assistance of others, with the exception of provided consultations. Within the master's thesis, the author must state the names of consultants and include a list of references.

Date of assignment receipt

Student's signature

Acknowledgements

At this place I would like to give a big thank you to supervisor of this master thesis, Karel Fliegel. His calm nature and expertise in field results in very valuable points and advice, during origin of this paper.

Declaration

Prohlašuji, že jsem předloženou práci vypracoval samostatně a že jsem uvedl veškeré použité informační zdroje v souladu s Metodickým pokynem o dodržování etických principů při přípravě vysokoškolských závěrečných prací.

I hereby declare that I have written this master thesis independently and quoted all the sources of information used in accordance with methodological instructions on ethical principles for writing an academic thesis. Moreover, I state that this thesis has neither been submitted nor accepted for any other degree.

Abstract

This paper aims to create a guideline in growing field of neural networks, and their application in learning based system for effective compression methods. In first chapter looking to create theoretical foundation, as to understand the terms that we will further work with - such as machine learning, (deep) neural networks, deep learning etc. Further we then discuss application of these methods in field of image coding as well as discussing criteria for classification of such systems. Second chapter, the practical part then discovers various implementation of learning based codecs (e.g. using TensorFlow platform) and their performance, comparing with conventional methods (JPEG, JPEG2000). The paper is wrapped up by discussing the results and future possibilities for research.

Keywords: machine learning, image processing, compression methods, neural networks

Supervisor: Ing. Fliegel Karel Ph.D.
Praha, Technická 2, B3-556

Abstrakt

Tato práce si klade za cíl vytvořit orientační návod stále rostoucím oborem strojového učení, neuronových sítí a jejich aplikace v systémech efektivních kompresních metod založených na učení. První kapitola tvoří teoretický základ, napomáhající k pochopení termínů se kterými se budeme nadále setkávat - t.j. strojové učení, (hluboké) neuronové sítě, hluboké učení atd. Dále budou diskutovány aplikace těchto metod právě v poli kódování obrazu, jakožto i návrh a diskuze kritérií pro klasifikaci těchto systémů. Následující kapitola, praktická část se věnuje různým implementacím kodeků, systémech založených na učení (např. s použitím platformy TensorFlow) a jejich efektivita v porovnání s konvenčními kompresními metodami (JPEG, JPEG2000). Práce je uzavřena diskuzí výsledků, návrhů zlepšení a možnostmi navazujícího výzkumu.

Klíčová slova: strojové učení, zpracování obrazu, kompresní metody, neuronové sítě

Překlad názvu: Metody strojového učení pro efektivní kompresi obrazu

Contents

| | |
|---|-----------|
| 1 Introduction | 1 |
| 2 Theoretical part | 3 |
| 2.1 Machine learning | 4 |
| 2.1.1 General idea | 4 |
| 2.1.2 Problem classification in Machine learning | 5 |
| 2.2 Artificial Neural network | 7 |
| 2.2.1 Biological background of Neural Networks | 8 |
| 2.2.2 Artificial neuron | 9 |
| 2.2.3 Learning and efficiency | 12 |
| 2.2.4 Deep Learning | 14 |
| 2.2.5 Topology of ANN | 14 |
| 2.3 Applications in image compression | 19 |
| 2.3.1 Neural networks and image coding | 20 |
| 2.3.2 Examples from practice | 22 |
| 2.3.3 Quality assessment - metrics . | 23 |
| 2.3.4 Data sets | 24 |
| 2.4 Classification proposal | 25 |
| 2.4.1 Topology | 25 |
| 2.4.2 Training | 26 |
| 2.4.3 Performance | 27 |
| 3 Practical part | 29 |
| 3.0.1 Models | 29 |
| 3.0.2 Metrics | 30 |
| 3.0.3 Performance evaluation | 31 |
| 4 Discussion and conclusions | 41 |
| A List of electronic attachments | 43 |
| B Average objective scores for Kodak dataset | 45 |
| Bibliography | 47 |

Figures

| | |
|--|--|
| <p>2.1 A Venn diagram showing the relations of fields like AI, machine learning, deep learning etc. Adapted from [18] 3</p> <p>2.2 Examples of Supervised Learning (Linear Regression) and Unsupervised Learning (Clustering). Adapted from [43] 6</p> <p>2.3 Illustration of the brain. Adapted from [27] 8</p> <p>2.4 Simple illustration of perceptron. Adapted from [39] 9</p> <p>2.5 Information processed in neuron. Adapted from [27] 10</p> <p>2.6 Example of various activation functions. Adapted from [39] 11</p> <p>2.7 Illustration of multilayered neural network. Adapted from [39] 12</p> <p>2.8 3-D plot of two variable cost function. Adapted from [39] 13</p> <p>2.9 Feed forward neural network. Adapted from [57] 15</p> <p>2.10 Recurrent neural network. Adapted from [57] 16</p> <p>2.11 Node map of an autoencoder network. Adapted from [57] 17</p> <p>2.12 Sparse autoencoder map. Adapted from [57] 17</p> <p>2.13 Architecture of LeNet-5, a CNN, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical. Adapted from [29] 19</p> <p>2.14 Image generated by deconvolutional layers. Adapted from [24] 20</p> <p>2.15 Simplified diagram of JPEG compression. 20</p> <p>2.16 Example of JPEG compression. (Left) raw, uncompressed image, (middle) medium level compression, perceptually lossless, (right) highly compressed image. Adapted from [40] 21</p> | <p>2.17 Pretraining consists of stacking restricted Boltzmann machines (RBMs), each with one layer of feature detectors. After pretraining, the RBMs are "unrolled" to create deeper autoencoder, which is tuned using backpropagation of error derivatives. Adapted from [20]..... 22</p> <p>2.18 Single iteration of shared RNN architecture from [54] 23</p> <p>3.1 Visual comparison of reconstructed image Kodim19 from Kodak, bpp ≈ 0.11 31</p> <p>3.2 Area of interest (fence), cropped from reconstructed image Kodim19 from Kodak, bpp ≈ 0.11 32</p> <p>3.3 Visualization of block pattern artifact from HyperMSE model. .. 33</p> <p>3.4 Visual comparison of reconstructed image Kodim20 from Kodak, bpp ≈ 0.11 35</p> <p>3.5 Area of interest (logo), cropped from reconstructed image Kodim20 from Kodak, bpp ≈ 0.11 36</p> <p>3.6 Measured scores for various objective quality metrics (Kodim19). 37</p> <p>3.7 Measured scores for various objective quality metrics (Kodim19). 38</p> <p>3.8 Measured scores for various objective quality metrics (Kodim20). 39</p> <p>3.9 Measured scores for various objective quality metrics (Kodim20). 40</p> <p>B.1 Average objective scores over Kodak data set. 45</p> <p>B.2 Average objective scores over Kodak data set. 46</p> |
|--|--|



Chapter 1

Introduction

With gradually increasing requirements for effective representation and compression of audiovisual content, for purpose of storing and sharing, we are always looking for new ways and methods to push the boundary of target quality and compression ratio achievable. With help of learning based algorithms, in-depth research of neural networks and development of efficient algorithms to train them - this appears to be a promising field to be applied in new compression systems; this is also supported by JPEG Committee forming a group, focused on exploring learning based image coding, so called JPEG-AI, which is aiming to develop an end-to-end compression system with potential for a standard to be defined. In this paper we would like to discuss state of the art compression methods, their design, architecture, as well as their performance in comparison with conventional systems and methods that are being commonly used on daily basis. Further a classification criteria would also be discussed.

Chapter 2

Theoretical part

Before delving deep into the topic at hand, let us briefly mention some important terms that might often be intertwined or substituted - such as *Artificial Intelligence*, *Neural Networks*, *Deep Learning* and more. For better understanding, let us refer to Venn diagram 2.1 suggested by [18], that is depicting the possible mutual connections of these terms, that we will further be working with.

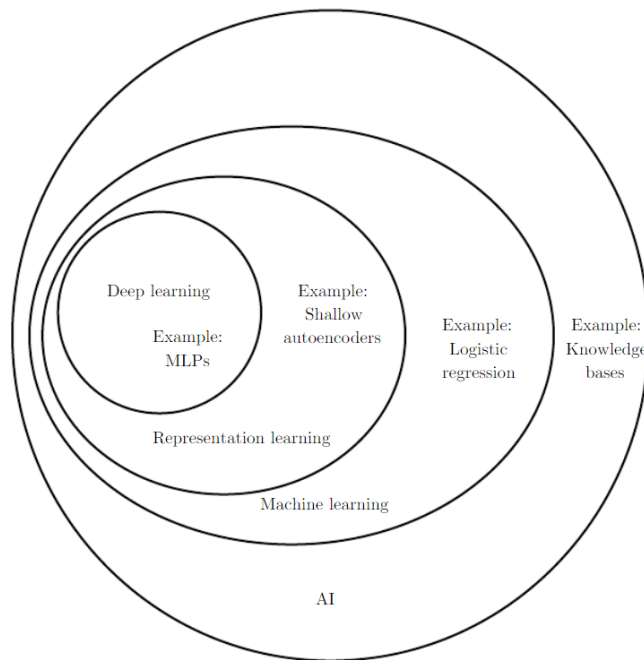


Figure 2.1: A Venn diagram showing the relations of fields like AI, machine learning, deep learning etc. Adapted from [18]

As is apparent from diagram above, the most superior field that we define and is related to this thesis (although way too broad) is **Artificial Intelligence** (AI). One of possible definitions define the field as the study of "intelligent agents", meaning any device that perceives its environment and takes actions that maximize its chance of successfully achieving its goals [42].

Formally, **machine learning** (ML) is a sub-field of artificial intelligence. However, in recent years, some organizations have begun using the terms artificial intelligence and machine learning interchangeably (hence the ambiguity of the terms). By ML we usually refer to a program or system that builds (trains) a predictive model from input data. The system uses the learned model to make useful predictions from new (never-before-seen) data drawn from the same distribution as the one used to train the model. It also can refer to the field of study concerned with these programs or systems [32]. We will discuss the topic in further detail in upcoming section 2.1.

Last two terms that we would like to mention in this section, which are closely linked is **Artificial Neural Networks** and **Deep Learning**. Former term refers to the biologically-inspired programming paradigm which enables a computer to learn from observational data. Latter term then describes a powerful set of techniques for learning in neural networks. Neural networks and deep learning currently provide the best solutions to many problems in image recognition, speech recognition, and natural language processing [39], and therefore being fundamental concepts for the aims of this thesis. These will be deeply mentioned in section 2.2 as well.

2.1 Machine learning

As our understanding of computers continues to mature, it seems inevitable that machine learning will play an increasingly central role in computer science and computer technology. [36] Let us describe some general ideas behind the term **Machine learning** (abbr. ML) and why is it perspective discipline in these mentioned fields.

2.1.1 General idea

With growing complexity of tasks that we deal with in modern computer science, programming, we naturally wish to use optimal approaches to reach the best possible efficiency. By using Machine learning methods, we gain time saving capabilities to process larger amount of data than ever before. Furthermore with traditional programming, our program might be limited or rather tailor-made to specific task that is defined apriori, without further adaptability to possible changes or exceptions that might come up in future. Machine learning is an application of artificial intelligence (AI) that provides systems the ability to automatically learn and improve from experience without being explicitly programmed. Machine learning focuses on the development of computer programs that can access data and use it learn for themselves. [60] This is to some extent analogous to human way of learning - we change our behaviour and approach to different tasks based on experience gained through our life.

To sum this up, let's use definition from [36]:

Definition 2.1.1 (Machine Learning). A computer program is said to learn from experience E with respect to some class of tasks T and performance

measure P , if its performance at tasks in T , as measured by P , improves with experience E .

As was mentioned before, this might not be the only possible definition, but for our purpose of better understanding and it works well enough.

Let us propose one example from practice. The most significant early milestone was A. L. Samuel's study using the game of checkers. Samuel devised detailed procedures both of "rote-learning" and "learning by generalization". When coupled with efficient methods of look ahead and search, these procedures enabled the computer to raise itself by prolonged practice from the status of a beginner to that of a tournament player. Hence there now exists a checkers program which can learn through experience of checkers to play better checkers. [34] In general, to have a well-defined learning problem, we must identify these three features: *the class of tasks, the measure of performance to be improved, and the source of experience*. To demonstrate this on our checkers learning problem:

- Task T : playing checkers
- Performance measure P : percent of games won against opponents
- Training experience E : playing practice games against itself

We can specify many learning problems in this fashion, such as learning to recognize handwritten words, or learning to drive a robotic automobile autonomously. [36]

■ 2.1.2 Problem classification in Machine learning

In section 2.1.1 we discussed general ideas behind the the ML - as well as some basics terms - with help of an example of simple checkers "program". Let us take a step back now, and generalize the problems (as in types of tasks we deal with) that we define in modern ML methods.

■ Supervised vs. Unsupervised Learning

Often, we talk about ML as having two paradigms - *supervised* and *unsupervised learning*. However, it is more accurate to describe ML problems as falling along a spectrum of supervision between supervised and unsupervised learning. Nevertheless, for the sake of simplicity, we will be describing these two extremes.

In **supervised machine learning**, you feed the features and their corresponding labels into an algorithm in a process called training. During training, the algorithm gradually determines the relationship between features and their corresponding labels. This relationship is called the model. Often times in machine learning, the model is very complex. Further it finds patterns between data and labels that can be expressed mathematically as functions. Given an input feature, you are telling the system what the expected output label is, thus you are supervising the training. The ML system will learn

patterns on this labeled data. In the future, the ML system will use these patterns to make predictions on data that it did not see during training [9]. Such example of ML can be seen on figure 2.2 on the left. Below we will present some of the most representative classifiers [43]:

- *Perceptron* and *Logistic Regression (LR)* are probably the simplest linear classifiers. For both models, the model (i.e., weights and bias) is basically a simple linear transformation.
- *Artificial Neural Networks (ANN)* is a general extension of the aforementioned linear classifiers. Compared with Perceptron or LR which linearly project input data to the output, ANN has an additional “hidden layer” (with a non-linear activation function), which enables ANN to model non-linearity. We will go deeper into this type of classifier in section 2.2.
- *Decision Tree (DT)* [44] and *Random Forest (RF)* [7] are two tree-structure based non-linear classifiers. Based on certain attribute-splitting criteria (e.g., Information Gain or Gini Impurity), DT can analyse the most informative attributes sequentially (i.e., splitting) until the final decision can be made.
- *Support Vector Machine (SVM)* [10] is another popular supervised learning method, also called large margin classifier, as it aims at finding a hyperplane that is capable of separating the data points (belonging to different classes) with the largest margin. For non-linearly separable data sets, various kernels (e.g., RBF (Radial Basis Function)) can be applied into the SVM framework with good generalization ability.
- In comparison to SVM, we have *K-Nearest Neighbour (KNN)* [11], which does not require a training process (also referred to as lazy learning), is another powerful non-linear classifier. The classification is performed by distance calculation (between query and all the training examples), distance ranking, and majority voting among the (K) nearest neighbours.

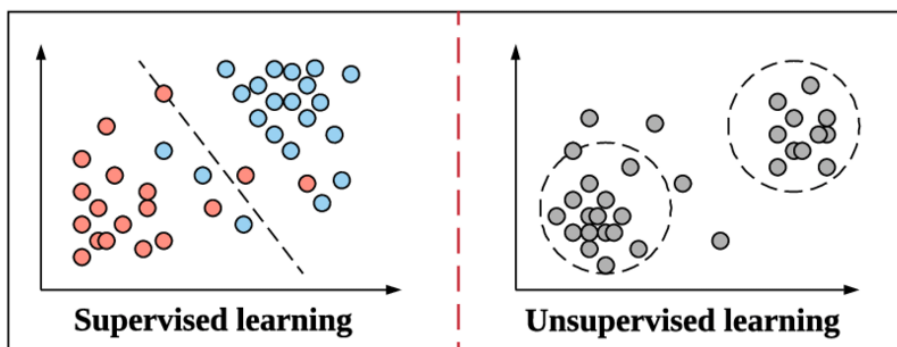


Figure 2.2: Examples of Supervised Learning (Linear Regression) and Unsupervised Learning (Clustering). Adapted from [43]

In **unsupervised learning**, the goal is to identify meaningful patterns in the data (refer to the figure 2.2, right). To accomplish this, the machine must learn from an unlabeled data set. In other words, the model has no hints how to categorize each piece of data and must infer its own rules for doing so [9]. For example, the clustering algorithm can be used to find the potential patterns of some unlabelled data and the obtained results can be used for future analysis. [43]

K-Means [19] and *Principal Component Analysis (PCA)* [49] are the two most popular unsupervised learning algorithms. K-means aims to find K group patterns from data by iteratively assigning each sample to different clusters based on the distance between the sample and the centroid of each cluster. PCA is normally used for dimensionality reduction, which can de-correlate the raw features before selecting the most informative ones.

■ Reinforcement Learning

Another quite specific approach that is also worth to mention is **reinforcement learning**. The term by itself is somewhat ambiguous, in that it refers to simultaneously a problem, a class of solution methods that work well on the class of problems, and the field that studies these problems and their solution methods.

Reinforcement learning problems involve learning what to do—how to map situations to actions—so as to maximize a numerical reward signal. In an essential way they are closed-loop problems because the learning system's actions influence its later inputs. Moreover, the learner is not told which actions to take, as in many forms of machine learning, but instead must discover which actions yield the most reward by trying them out [50]. The lack of a data requirement makes RL a tempting approach. However, designing a good reward function is difficult, and RL models are less stable and predictable than supervised approaches [9].

To put this simply reinforcement learning focuses on "trial and error" style of training, by discovering best possible ways to reach the predefined goal of our task, while receiving positive or negative feedback, to determine how well is the learning agent doing.

■ 2.2 Artificial Neural network

The Artificial Neural Networks (ANN) or more often simply **Neural Networks** are computing systems vaguely inspired by the biological neural networks [8]. Before we begin with the technical description, it would be useful to start with brief description of the biology of neural networks as we know it from living organisms.

2.2.1 Biological background of Neural Networks

The entire information processing system, i.e. the vertebrate nervous system, consists of the central nervous system and the peripheral nervous system, which is only a first and simple subdivision. In reality, such a rigid subdivision does not make sense, but here it is helpful to outline the information processing in a body.

The *peripheral nervous system (PNS)* comprises the nerves that are situated outside of the brain or the spinal cord. These nerves form a branched and very dense network throughout the whole body.

We would naturally like to focus on brain, as the main center of information processing. To keep things simple, we will divide brain into 4 main parts, as seen on 2.3.

The *cerebrum* (telencephalon) is one of the areas of the brain that changed most during evolution. Along an axis, running from the lateral face to the back of the head, this area is divided into two hemispheres, which are organized in a folded structure.

These cerebral hemispheres are connected by one strong nerve cord ("bar") and several small ones. A large number of **neurons** are located in the *cerebral cortex* (cortex) which is approx. 2-4 cm thick and divided into different *cortical fields*, each having a specific task to fulfill. *Primary cortical fields* are responsible for processing qualitative information, such as the management of different perceptions (e.g. the *visual cortex* is responsible for the management of vision). *Association cortical fields*, however, perform more abstract association and thinking processes; they also contain our memory. [27]

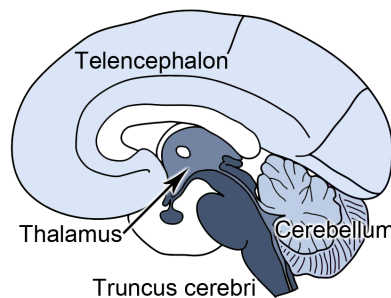


Figure 2.3: Illustration of the brain. Adapted from [27]

To close out this section, we will focus on fundamental description of said neuron, which we can consider as a basic processing unit the brain. A neuron is nothing more than a switch with information input and output. The switch will be activated if there are enough stimuli of other neurons hitting the information input. Then, at the information output, a pulse is sent to, for example, other neurons. To put it simply:

- *Synapses* weight the individual parts of information
- *Dendrites* then collect all parts of this information

- In the *soma* the weighted information is accumulated
- And the axon transfers outgoing pulses

All of the parts mentioned above can be found depicted in figure

2.2.2 Artificial neuron

Now let us put the knowledge gained in previous section 2.2.1, into context of basic concepts of ANNs, their structure and design.

The fundamental unit the forms our neural network is once again neuron, be it artificial one, called **perceptron**. To be more precise, perceptron is a simple artificial neuron whose *activation function* consists of taking the total net input and outputting 1 if this is above a threshold T , and 0 otherwise[61]. We will describe what is activation function further in this chapter. Perceptrons were developed in the 1950s and 1960s by the scientist Frank Rosenblatt, inspired by earlier work by Warren McCulloch and Walter Pitts. Today, it's more common to use other models of artificial neurons, and the main neuron model used is one called the *sigmoid neuron* [39].

Very similarly as with biological neuron, a perceptron takes several binary inputs, x_1, x_2, \dots, x_n and produces a single binary output:

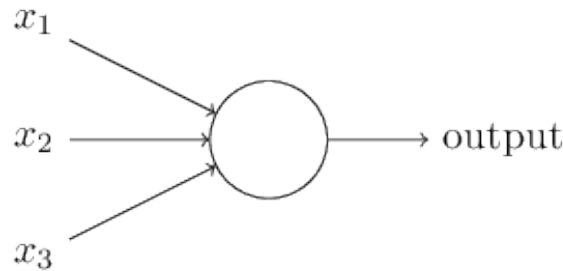


Figure 2.4: Simple illustration of perceptron. Adapted from [39]

Example above on figure 2.4 shows perceptron with three inputs x_1, x_2, x_3 , which is in this case completely arbitrary, and in general perceptron can have fewer or more inputs. Rosenblatt proposed a simple rule to compute the output. He introduced weights, w_1, w_2, \dots, w_n real numbers expressing the importance of the respective inputs to the output. The neuron's output, 0 or 1, is determined by whether the weighted sum $\sum_j w_j x_j$ is less than or greater than some threshold value, arithmetically put:

$$\text{output} = \begin{cases} 0, & \text{if } \sum_j w_j x_j \leq \text{threshold} \\ 1, & \text{if } \sum_j w_j x_j > \text{threshold} \end{cases} \quad (2.1)$$

That's the basic mathematical model, although perceptron isn't a complete model of human decision-making. Connection between the neurons carries the information that is to be processed. Naturally our next question should be - how is the incoming information processed? For our neuron j , [27] defines *propagation function* converting vector inputs to scalar network inputs:

Definition 2.2.1 (Propagation function and network input). Let $I = \{i_1, i_2, \dots, i_n\}$ be the set of neurons, such that $\forall z \in \{1, \dots, n\} : \exists w_{i_z, j}$. Then the network input of j , called net_j , is calculated by the propagation function f_{prop} as follows:

$$\text{net}_j = f_{\text{prop}}(o_{i_1}, \dots, o_{i_n}, w_{i_1, j}, \dots, w_{i_n, j}) \quad (2.2)$$

Here the *weighted sum* is very popular: The multiplication of the output of each neuron i by $w_{i, j}$, and the summation of the results:

$$\text{net}_j = \sum_{i \in I} (o_i \cdot w_{i, j}) \quad (2.3)$$

At every moment, every neuron in network happens to be in some state, it's "active" so to speak. This state further defines its reaction on incoming information, and we refer to it as activation state or more often just shortly *activation*. [27] provides following general definition:

Definition 2.2.2 (Activation state / activation). Let j be a neuron. The activation state a_j is explicitly assigned to j , indicates the extent of the neuron's activity and results from the activation function.

So the current activation state of neuron is based on the activation state as well as the input, incoming information. This relation is described by *activation function* [27]:

Definition 2.2.3 (Activation function). Let j be a neuron. The activation function is defined as

$$a_j(t) = f_{\text{act}}(\text{net}_j(t), a_j(t-1), \Theta_j). \quad (2.4)$$

It transforms the network input net_j , as well as the previous activation state $a_j(t-1)$ into a new activation state $a_j(t)$, with the threshold value playing an important role.

Unlike the other variables within the neural network the activation function is often defined globally for all neurons or at least for a set of neurons and only the threshold values are different for each neuron. The threshold value can also be changed during time, e.g. by learning procedure. [27] As for the most common activation functions: the most simple one is the *binary threshold function*, also well known as Heaviside function (refer to figure 2.6) - value changes from one to another, once overcoming given threshold; otherwise it stays constant. Other very popular functions are *Fermi function* (also logistic function), possibly expanded by a temperature agent T

$$\frac{1}{1 + e^{-\frac{x}{T}}} \quad (2.5)$$

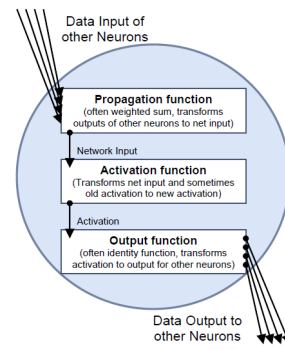


Figure 2.5: Information processed in neuron. Adapted from [27]

which gives us the possibility to change the gradient of the function and approximate the Heaviside function. Another function would be *hyperbolic tangent*. Both lastly mentioned, in contrast to Heaviside function, are differentiable (fig. 2.6). Incidentally, there exist activation functions which are not explicitly defined but depend on the input according to a random distribution - *stochastic activation function*. [27]

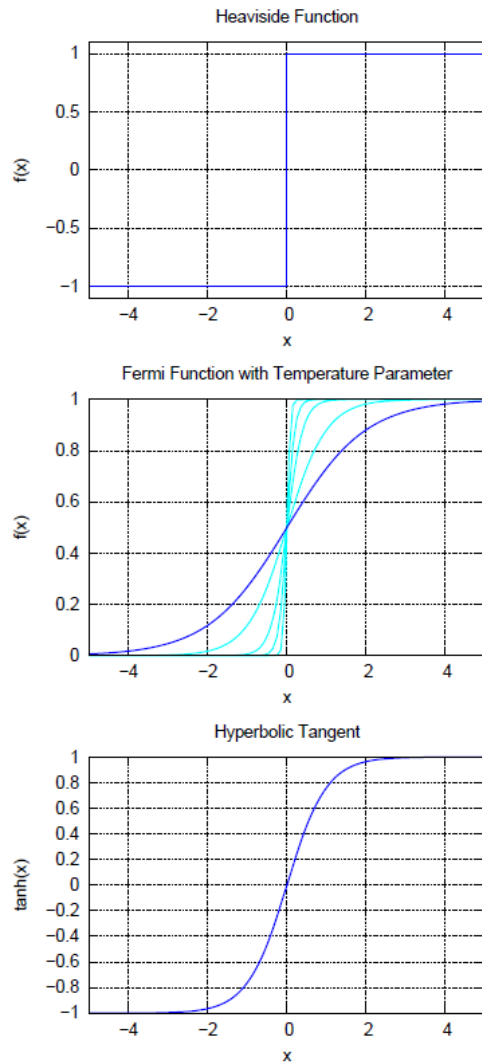


Figure 2.6: Example of various activation functions. Adapted from [39]

Lastly we define *output function* of a neuron j calculates the values which are transferred to the other neurons connected to j :

Definition 2.2.4 (Output function). Let the j be a neuron. The output function

$$f_{\text{out}}(a_j) = o_j \tag{2.6}$$

calculates the output value o_j of the neuron j from it's activation state a_j .

Generally, the output function is defined globally, too. Often this function is the identity, i.e. the activation a_j is directly output:

$$f_{\text{out}}(a_j) = a_j, \text{ so } o_j = a_j \quad (2.7)$$

Let us illustrate an example, how a perceptron can weigh up different kinds of evidence in order to make decisions - this is plausible with a complex network of perceptrons that could make quite subtle decisions:

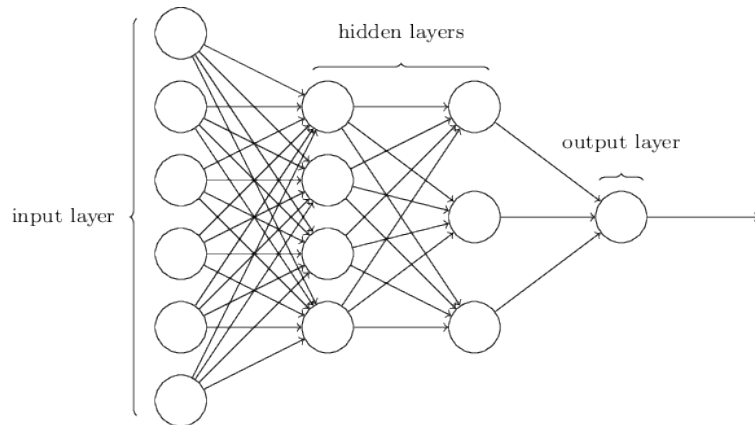


Figure 2.7: Illustration of multilayered neural network. Adapted from [39]

The leftmost layer in this network is called the input layer, and the neurons within the layer are called input neurons. The rightmost or output layer contains the output neurons, or, as in this case, a single output neuron. The middle layer is called a hidden layer, since the neurons in this layer are neither inputs nor outputs [39]. Therefore an example in figure 2.7 is four-layer network, with two hidden layers.

While the design of the input and output layers of a neural network is often straightforward, there can be quite an art to the design of the hidden layers. In particular, it's not possible to sum up the design process for the hidden layers with a few simple rules of thumb. Instead, neural networks researchers have developed many design heuristics for the hidden layers, which help people get the behaviour they want out of their nets. For example, such heuristics can be used to help determine how to trade off the number of hidden layers against the time required to train the network. [39]

■ 2.2.3 Learning and efficiency

Since we are aiming to create an efficient model suitable to given task, we'd like is an algorithm which lets us find *weights* and *biases* so that the output from the network approximates $y(x)$ for all training inputs x . To quantify how well we're achieving this goal we define a **cost function**¹:

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a\|^2 \quad (2.8)$$

¹Sometimes referred to as a loss or objective function.

Here, w denotes the collection of all weights in the network, b all the biases, n is the total number of training inputs, a is the vector of outputs from the network when x is input, and the sum is over all training inputs, x . This is actually nothing else than *mean squared error (MSE)*². This means that objectively the goal of our training algorithm would be to keep the cost function as low as possible, $C(w, b) \approx 0$. The lower the MSE value, the closer the output value is to the correct result.

The next milestone on our way is - how do we find this minimum of cost function? This is reached by using an algorithm called *gradient descent*. Suppose we are trying to minimize cost function $C(v)$. This could be any real-valued function with arbitrary amount of values: let's say C is a function of m variables, v_1, v_2, \dots, v_m . Then change ΔC in C produced by a small change $\Delta v = (\Delta v_1, \dots, \Delta v_m)^T$ is

$$\Delta C \approx \nabla C \cdot \Delta v, \tag{2.9}$$

where gradient ∇C is the vector

$$\nabla C \equiv \left(\frac{\partial C}{\partial v_1}, \dots, \frac{\partial C}{\partial v_2} \right)^T. \tag{2.10}$$

We will further also define

$$\Delta v = -\eta \nabla C, \tag{2.11}$$

where η is a small, positive parameter (known as the learning rate). Then equation 2.9 tells us that $\Delta C \approx -\eta \nabla C \cdot \nabla C = -\eta \|\nabla C\|^2$. Because $\|\nabla C\|^2 \geq 0$, this guarantees that $\Delta C \leq 0$, i.e., C will always decrease, if we change v according to the prescription in 2.11.

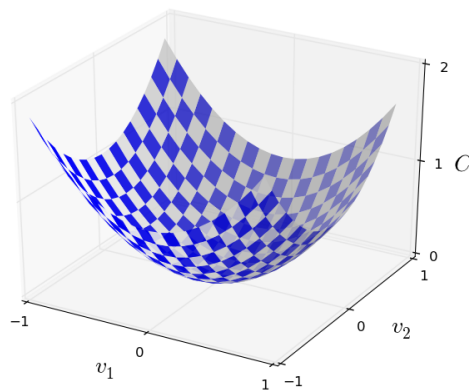


Figure 2.8: 3-D plot of two variable cost function. Adapted from [39]

²Or quadratic cost function.

To put this into a better perspective, let's present a case, where $C(v)$ being function of two variables v_1, v_2 and as we know, our goal is to find a global minimum of said function $C(v)$. By one glance at 3-D plot 2.8 of presented function $C(v)$, we can immediately recognize there exist ideal combination of variables v_1, v_2 , the would give us the desired result [39].

The way the gradient descent algorithm works is - by making small changes (defined by step η) and evaluating change Δv using equation 2.11 - to repeatedly compute the gradient ∇C , and then to move in the opposite direction, "falling down" the slope of the valley. This concept is then upscaled to arbitrary amount of values, which takes us back to 2.10.

■ 2.2.4 Deep Learning

It is worth mention what do we understand by the term **deep learning** and deep neural networks (DNNs) since this term is being more and more common while talking about modern learning-based methods.

Since idea of ANNs, that we started to described in previous sections, isn't anything new, it's no surprise that the whole field developed during the years. More often then not, we can see neural networks that doesn't have just a few stages, let's say single hidden layer - that is something we call *shallow neural network* models. We often look for more complex solutions to describe our task; extracting smaller, more detailed (but also more abstract) features. Lately we often achieve this by stacking hidden layers, leading to so called *deep neural networks* (DNNs).

Unsurprisingly this also brings a question of how to (effectively) train such NNs. Learning (or sometimes called credit assignment) is about finding weights (as described in section 2.2.3) that make the NN exhibit desired behavior, such as driving a car for instance. Depending on the problem and how the neurons are connected, such behavior may require long causal chains of computational stages, where each stage transforms (of ten in a non-linear way) the aggregate activation of the network. Deep Learning is about accurately assigning credit across many such stages. BP-based training of deep NNs with many layers, however, had been found to be difficult in practice by the late 1980s and had become an explicit research subject by the early 1990s. DL became practically feasible to some extent through the help of Unsupervised Learning (section 2.1.2), further in he 1990s and 2000s also saw brought many improvements of purely supervised DL and lastly Deep NNs also have become relevant for the more general field of Reinforcement Learning (2.1.2). [46]

■ 2.2.5 Topology of ANN

This section is going to focus on common topologies of neural networks and their names, often abbreviated. The fact that abbreviations are so commonly spread in this field, makes it rather hard and overwhelming for inexperienced person to orientate himself. Following chapter will often refer to very useful sheet by Fjodor van Veen [57], providing well sorted information and visual

representation of topic at hand. Another note - list that will follow further below won't (and almost can't be) complete, but should provide reader with basic idea of structure and design of NNs, especially those applicable for image processing (compression).

■ Feed Forward Neural Networks

Feed forward neural networks (FF or FFNN) and perceptrons represent the simplest examples of neural networks. We already described perceptrons in more details under section 2.2.2, but very briefly: it is the fundamental unit of neural networks, defined by weighted inputs from other neurons, activation function and its output.

Putting neurons into layers, we create network. Each layer containing input, hidden or output cells in parallel. FF network, as the name suggests, is a one-directional network - feeding the information from front and it propagates to the end, towards output cells on the right. All nodes are fully connected, meaning that every node in given layer has connections to all nodes in the next one. The applications of FFNNs are somewhat limited and they are often paired with other types of networks for more complex tasks. By itself, they can be used as model for classification.

The decision-making of the network lies in the hidden layer in the middle. As was mentioned in 2.2.2, every neuron has weights, which represents the decision-making capabilities, once the information is being propagated through them. These weights are set by process called *training*, during which we feed the network the data on the input and pair them with expected data on the output (so called supervised learning, 2.1.2). We can go back, to the input layer, layer by layer comparing layer input with layer error. The error being propagated is often some variation of the difference between the input and the output (like MSE or just the linear difference) [57]. Each layer, no matter how deep it is, contributes to next layer error; therefore, we can adjust weights by value we can get by multiplying previous layer output by output error multiplied by current layer output.[51] This process is called **Backpropagation**.

Amongst some variations to general FF network we put: **Radial Basis Network** (RBF), which is FF network that uses *radial basis function*³ as an activation function.2.2.3.

Another example is **Deep feed forward network** (DFF) pioneering the way for deep learning in early 90s. General idea is the same as for FFNNs, adding more hidden layers (therefore *deep*). Although stacking hidden layers

Feed Forward (FF)

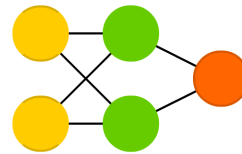


Figure 2.9: Feed forward neural network. Adapted from [57]

³real-valued function φ whose value depends only on the distance between the input and some fixed point, either the origin, or some other fixed point \mathbf{c} [13]

led to exponential growth of training time causing them to be impractical, with forthcoming more efficient approaches in early 00s, they now form a core of modern ML systems.[51][52]

■ Recurrent Neural Networks

Recurrent Neural Networks (RNN) brings new type of nodes, cells - recurrent cells. Our goal is to represent time implicitly by its effects on processing rather than explicitly (as in a spatial representation) [15]. This idea first described by Jordan, M.I. in 1986 [23] allows us to involve the use of recurrent links in order to provide networks with a dynamic memory.

Neurons are fed information not just from the previous layer but also from themselves from the previous pass. This also means that order in which data are fed to the network matter - feeding it "dog" then "cat" may yield different results compared to opposite order. One big problem with RNNs is the vanishing (or exploding) gradient problem where, depending on the activation functions used, information rapidly gets lost over time, just like very deep FFNNs lose information in depth. To put it in perspective: if the weight reaches a value of 0 (vanishing) or 1 000 000 (exploding), the previous state of the cell won't be very informative. In general, recurrent networks are a good choice for advancing or completing information, such as autocompletion. [57]

Variation of RNNs called **Long / short term memory (LSTM)** aims to deal with mentioned problem with exploding/vanishing gradient. This is achieved by an efficient, gradient-based algorithm for an architecture enforcing constant (thus neither exploding nor vanishing) error flow through internal states of special units - *memory cells*. [21] These cells are made composed of so called *gates*, which define whether to feed the information forward (and how much of it) or erase/forget it. [57] They are used in fields such as speech and writing recognition.

Gated recurrent units (GRU) then are very much similar to LSTMs by design. They are wired slightly differently, with different gating - introducing *update* and *reset* gates. GRUs are slightly less expressive, but a bit faster; commonly used in speech (or sound in general) synthesis.

To all networks mention above in this section also exist a bidirectional variation, called **Bidirectional recurrent neural networks, bidirectional long / short term memory networks and bidirectional gated recurrent units (BiRNN, BiLSTM and BiGRU respectively)**. As for the topology itself, they appear to be all them same. The difference is in fact that they can be trained without the limitation of using input information just up to a preset future frame. This is accomplished by training it simultaneously in

Recurrent Neural Network (RNN)

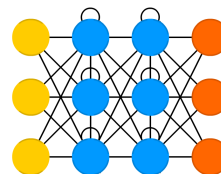


Figure 2.10: Recurrent neural network. Adapted from [57]

positive and negative time direction. This trains the network to fill in gaps instead of advancing information, e.g. instead of expanding an image on the edge, it could fill a hole in the middle of an image. [47][57]

■ **Autoencoders**

Auto Encoder (AE)

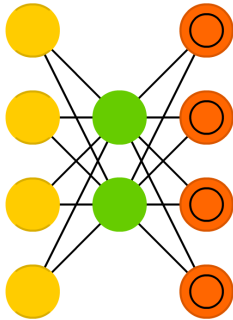


Figure 2.11: Node map of an autoencoder network. Adapted from [57]

Autoencoders (AE) are considerably similar by their design to FFNNs, applied a bit differently. Their aim to encode (compress) the information on the input. Structurally they resemble the hourglass-like shape (as can be seen on 2.11), and are symmetrical around the (hidden) layer in the middle. This smallest layer (or multiple layers potentially) are the place with the highest grade of compression of the information, and it is so called *bottleneck* or *chokepoint* of the network. Everything up to the middle layer is called the encoding part, everything after the decoding and the middle the code. [57] Bourlard (2000) shows that the nonlinearities of these hidden units are useless, and that the optimal parameter values can be derived directly by purely linear techniques relying on singular

value decomposition and low rank matrix approximation, similar in spirit to the well-known Karhunen-Loève transform (KLT). [6] Autoencoders are normally used for classification, clustering and feature compression.[52]

Denosing autoencoders (DAE) is variation of AE networks, with a little spin. Vincent et al. (2008) is investigation specific criterion of AEs called *robustness to partial destruction of the input* (basically ability to restore data with presence of noise on the input). The motivation behind is such: a good representation is expected to capture stable structures in the form of dependencies and regularities characteristic of the (unknown) distribution of its observed input. For high dimensional redundant input (such as images) at least, such structures are likely to depend on evidence gathered from a combination of many input dimensions. They should thus be recoverable from partial observation only. [58]

In simple terms this means that by training with the presence of noise, the model is encouraged to focus on more broad features in the data (e.g. image), rather than details, which are often shifting, changing. This leads to higher robustness to noise.

Sparse autoencoders (SAE) is in a way opposite to standard AEs. The idea is that instead of compressing the information in the middle (hidden) layer containing less nodes than input/output layer, we expand the information through bigger amount of nodes, as seen in 2.12. These types of networks

Sparse AE (SAE)

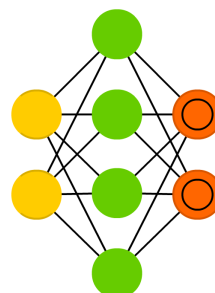


Figure 2.12: Sparse autoencoder map.

Simply put we start with *scanning layer* with a smaller dimension in comparison to input, that is moving (scanning) over the input image, feeding these data into convolutional layers. These convolutional layers also tend to shrink as they become deeper, mostly by easily divisible factors of the input (e.g. if we have layer of 20x20, it's likely that would be followed by 10, and 5). Besides these convolutional layers, they also often feature pooling layers. Pooling is a way to filter out details: a commonly found pooling technique is max pooling, where we take say 2 by 2 pixels and pass on the pixel with the most amount of red. [57]

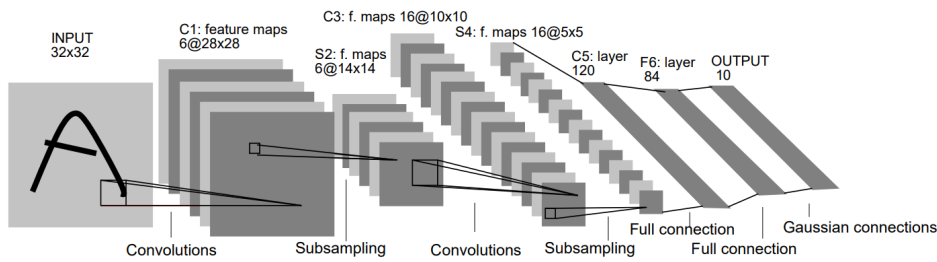


Figure 2.13: Architecture of LeNet-5, a CNN, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical. Adapted from [29]

Generative Adversarial Networks

Framework for estimating generative models via an adversarial process, in which we simultaneously train two models: a generative model G that captures the data distribution, and a discriminative model D that estimates the probability that a sample came from the training data rather than G . The training procedure for G is to maximize the probability of D making a mistake. This framework corresponds to a minimax two-player game. In the space of arbitrary functions G and D , a unique solution exists, with G recovering the training data distribution and D equal to $frac{1}{2}$ everywhere. In the case where G and D are defined by multilayer perceptrons, the entire system can be trained with *backpropagation*. [17]

In simple terms: GAN is actually two networks designed to compete against each other. Discriminative network D determines whether the input produced by the generative network G is good enough (difference from reference data); through backpropagation G adjust its weights until maximal quality solution.

2.3 Applications in image compression

With growing amount of visual data over the internet as well as the higher quality of this content (higher resolution such as UHD, high dynamic range etc.) the requirement of efficient compression and storing of this content becomes crucial. The research and development of image/video coding can be

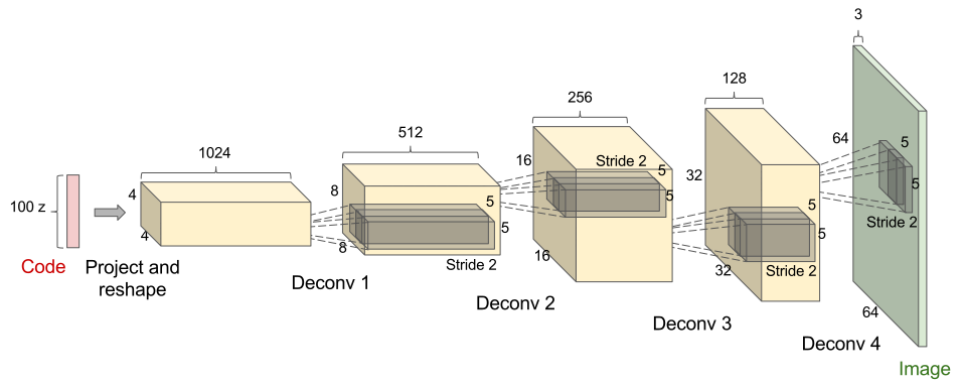


Figure 2.14: Image generated by deconvolutional layers. Adapted from [24]

dated back to as early as 1960s, much earlier than the appearance of modern imaging, image processing, and visual communication systems. [30]

Compression may or may not give us the possibility of (mathematically) perfect reconstruction; which is referred to as lossless or lossy compression. For natural image/video, the compression efficiency of lossless coding is usually below requirement, so most of efforts are devoted to lossy coding.

The redundancies within images and videos are fundamentally important for image and video compression, including spatial redundancy, visual redundancy and statistical redundancy. Besides, the temporal redundancy existing in video sequences enables the video compression to achieve higher compression ratio compared with image compression. [31]

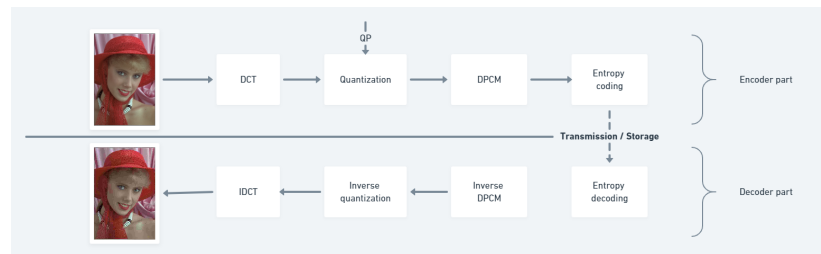


Figure 2.15: Simplified diagram of JPEG compression.

2.3.1 Neural networks and image coding

The past decade has witnessed the emerging and booming of deep learning, a class of techniques that are increasingly adopted in the hope of approaching the ultimate goal of artificial intelligence [28]. One of the most benefiting factors of using compression methods based on deep learning is that extracted representations of image/video content isn't manually designed, rather learned from visual data presented during training, learning procedure. Achieved high level of abstraction gives us possibility to extract latent representations of the image, as opposed to manually designed systems (e.g. quantization matrixes of JPEG systems). This is further useful especially for processing natively



Figure 2.16: Example of JPEG compression. (Left) raw, uncompressed image, (middle) medium level compression, perceptually lossless, (right) highly compressed image. Adapted from [40]

unstructured data, such as acoustic and visual signal, whilst processing such data has been a longstanding difficulty in the artificial intelligence field. [30]

In general we can divide available compression systems into two main groups suggested by Liu et al. in [30]:

- Deep schemes - new coding schemes, that are built primarily upon deep networks
- Deep tools - used within traditional coding schemes or together with traditional coding tools.

To put simply deep schemes can be viewed as an end-to-end compression system based on NNs and deep learning (pixel probability modeling and auto-encoders), whereas deep tools are singular tools designed to replace specific part of the traditional systems: intra-picture prediction, inter-picture prediction, cross-channel prediction, probability distribution prediction, transform, post- or in-loop filtering, down- and up-sampling, as well as encoding optimizations. [30]

Auto-encoder we briefly described in section 2.2.5. It originates from work of Hinton and Salakhutdinov [20], which explores how high-dimensional data can be converted to low-dimensional codes (encoding) by training a multilayer neural network with a small central layer to reconstruct high-dimensional input vectors (decoding).

When we consider the compression requirement, there are several challenges [30]:

- the low-dimension representation shall be quantized then coded, but the quantization step is not differentiable, making a difficulty to train the network.
- lossy coding is to achieve a better tradeoff between rate and quality, so the rate shall be taken into account when training the network, but the rate is not easy to calculate or estimate.

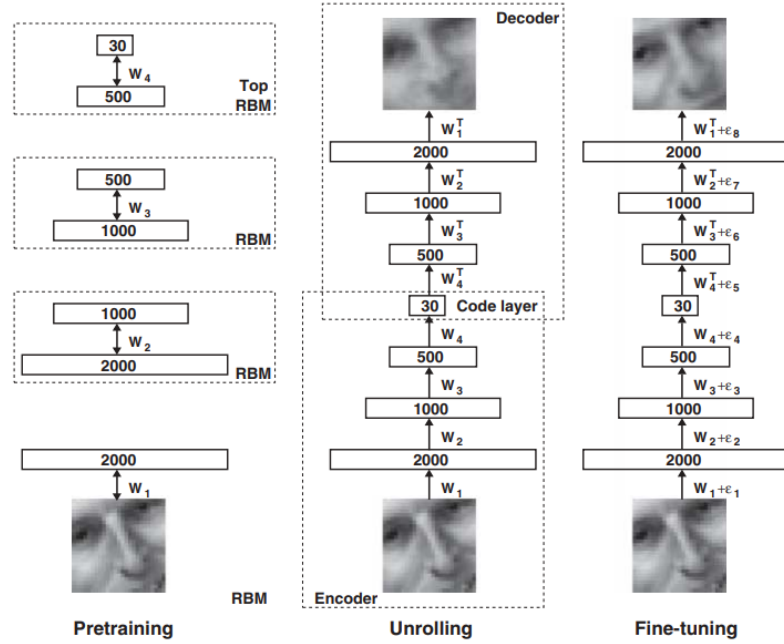


Figure 2.17: Pretraining consists of stacking restricted Boltzmann machines (RBMs), each with one layer of feature detectors. After pretraining, the RBMs are "unrolled" to create deeper autoencoder, which is tuned using backpropagation of error derivatives. Adapted from [20]

- a practical image coding scheme needs to consider variable rate, scalability, encoding/decoding speed, interoperability, and so on.

The original image x is transformed to $y = g_a(x)$, and y is quantized then coded. The decoded \hat{y} is inversely transformed to $\hat{x} = g_s(\hat{y})$. Considering the tradeoff between rate and quality, we can train the network to minimize the joint rate-distortion cost $D + \lambda R$ where D is calculated or estimated as the difference between x and \hat{x} (note that the difference may be calculated or estimated in a perception space), R is calculated or estimated from the quantized code, and λ is the Lagrange multiplier. All of the existing researches follow this scheme more or less and differ in their network structure and loss function. [30]

■ 2.3.2 Examples from practice

This section should represent some real implementations that have been proposed, and that would further be part of a performance comparison in practical part of this paper.

Paper of Minnen et al. (2018) [35] introduces they learning-based compression method; they network basically consist of two sub-networks, each having a different purpose. The first is the core *autoencoder*, which learns a quantized latent representation of images (Encoder and Decoder blocks). The second sub-network is responsible for learning a probabilistic model over

quantized latents used for entropy coding. It combines the Context Model, an autoregressive model over latents, with the hyper-network (Hyper Encoder and Hyper Decoder blocks), which learns to represent information useful for correcting the context-based predictions. The data from these two sources is combined by the Entropy Parameters network, which generates the mean and scale parameters for a conditional Gaussian entropy model.

Toderici et al. (2016) proposed an architecture (fig. 2.18) with RNN-based encoder and decoder, a binarizer, and a neural network for entropy coding. Their goal was to provide a neural network which is competitive across compression rates on images of arbitrary size. Also most of autoencoders of the time were defined by fixed compression rate, based on the size of bottleneck layer in the middle. [54] extends this idea by supporting variable rate compression while maintaining high compression rates beyond thumbnail-sized images.

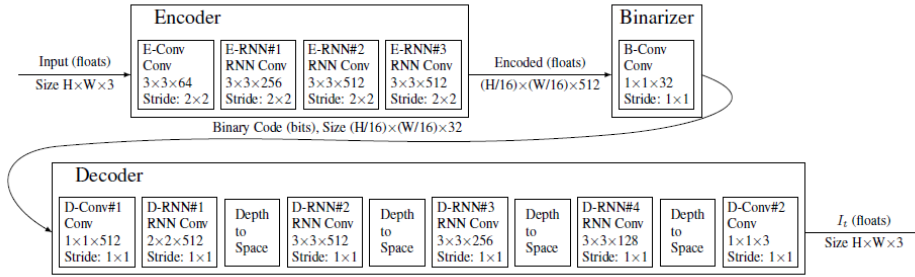


Figure 2.18: Single iteration of shared RNN architecture from [54]

While the network weights are shared between iterations, the states in the recurrent components are propagated to the next iteration. Therefore residuals are encoded and decoded in different contexts in different iterations. Proposed network incorporates LSTM units, Associative LSTMs⁵ (experimentally proved to be effective only in decoder) and GRU - computation block that passes residual information around the block in order to speed up convergence.

2.3.3 Quality assessment - metrics

The most widely used full-reference⁶ image quality and distortion assessment algorithms are peak signal-to-noise ratio (PSNR) and mean squared error (MSE), which do not correlate well with perceived quality. That is a main reason for development of method called **Multi-scale structural similarity index measure** (abbr. MS-SSIM), that will be further used as main metric of quality.

Structural similarity provides an alternative and complementary approach to the problem of image quality assessment.[64][63] The perceivability of

⁵Associative LSTM extends LSTM using holographic representation

⁶based on an initial uncompressed or distortion-free image as reference

image details depends the sampling density of the image signal, the distance from the image plane to the observer, and the perceptual capability of the observer's visual system. A older single-scale method may be appropriate only for specific settings and is surpassed by multi-scale version that is a convenient way to incorporate image details at different resolutions. [59]

As was mentioned before, more simple metrics such as **MSE** or **PSNR** might not as descriptive for the subjective quality assessment, nevertheless they still might be useful as indicators, e.g. during training our network. So for the image I ($m \times n$ pixels) and its reconstruction K , we define:

$$\text{MSE} = \frac{1}{mn} \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} [I(i, j) - K(i, j)]^2 \quad (2.12)$$

$$\text{PSNR} = 10 \cdot \log\left(\frac{\text{MAX}_I^2}{\text{MSE}}\right) \quad (2.13)$$

For the purpose of new learning based compression methods, these previously mentioned metrics might not be fully optimal. For this reason, new metrics for objective quality assessment are developed as well. Work of Testorina et al. [53] gives extensive comparison of objective metrics available and their correlation to subjective tests. Based on various correlation coefficients such as Pearson Linear Correlation Coefficient (PLCC), Spearman rank order correlation coefficient (SROCC) and Kendall's Rank Correlation Coefficient (KROCC) - few metrics stands out. Therefore let us extend our list.

The Visual Information Fidelity (VIF or VIFP) [48] express the quality of the image using Natural Image Statistics exploiting also the Human Visual System characteristics. The original code works with single-channel images. The final RGB score is computed as the average over the 3 channels. The metric takes values in range 0-1 and a high score express better image quality.

The Feature-Similarity Index Metric (FSIM) [62] uses the Phase Congruency (PC) and Gradient Magnitude (GM) to assess image local quality. Since this metric allows as input only monochromatic images, the score has been computed over the R, G and B channel separately as well as on the grayscale version of the original image. The FSIM for the RGB image is computed as the average of the one on the R, G and B results. A high metric value express better image quality.

In paper of Ponomarenko [41] a simple model of between-coefficient masking of DCT basis functions is proposed and the modifications of PSNR that take into account this model are put forward. The new measure **PSNR-HVS-M** has shown its higher efficiency (adequacy) in comparison to known metrics.

■ 2.3.4 Data sets

To precisely define the objective quality of compression method at hand, we should also use suitable testing data, broad enough, various so we can define what and where are the strengths and weaknesses.

From my research, the most common data set used by many papers is still **Kodak Lossless True Color Image Suite**[26], which contains 24 images in so called "full colour" (24 bits per pixel, 8 bits per colour layer) contained in lossless format PNG. Each image is of size 768×512 pixels (both dimensions divisible by 32, commonly used factor in compression systems). This set was released by the Kodak Corporation for unrestricted research usage in 1999.

Although possibly prevalent, it's hardly the only data set, so let us also propose other options.

- **Tecnick** [2] containing data set of 100 images, size 1200×1200 .
- **MCL-JCI** [22] has 50 images of 1920×1080 with 8-bit depth that can be classified into ten semantic categories such as people, animals, plants, buildings, water or lake, sky, bridge, transportation vehicles (boats or cars) and indoor.
- **LIVE in the Wild** [16] with 1162 distorted images from different mobile devices (therefore different sizes) Images of faces, people, animals, close-up shots, wide-angle shots, nature scenes, man-made objects, images with distinct foreground/background configurations and without any object of interest.

Many others could be listed, the choice heavily depends on the expected task, what would be the subject of our study.

■ 2.4 Classification proposal

This paper also aims to propose a taxonomy, how to classify and orient yourself in vast field of neural networks, based on various criteria. This would also help us with even field on which we could compare various networks and their implementations. We have to consider that we are discussing very broad field and different perspectives should be pointed out. Nelson and Rogers (1992) [38] points out, that there might be different criteria to choose as *best* per se, for the same task. Therefore the final criteria should be selected by the researcher or application specialist depending on the results desired. Their paper also proposes division of criteria into 3 main groups - topology, training and performance.

■ 2.4.1 Topology

Let us start with topological point of view; where we are considering the architecture, design of the network itself. In section 2.2.5 we've been describing different types of network topologies [14] simply defines category *Neural network type*, i.e. top level description of the network. Our aim is to create efficient compression method - do we use AE network? Or some variation, such as VAE? Or possibly we could train model through combination of networks, e.g. GANs? This is also connected to what type of units/neurons/nodes do

we define in our network and their activation functions (2.2.3). Also, what is the extent of our network in means of end-to-end image compression system - NN focused on encoding residual part, the entropy coder or combination of both?

If we were to dissect this idea into smaller problems, we should discuss the complexity of the network. **How many (hidden) layers** should our NN - shallow or deep NN? Considering solely the amount of layers, while assuming that same performance regardless of this amount, we would prefer *fewest layer amount possible*.

How many nodes do we have per layer? Generally speaking the more nodes a network has the greater the tendency to memorize the training set and the less tendency to properly generalize. The network with the *fewest nodes*, irrespective of how they are arranged or connected could be selected by this criterion. [38]

And lastly, what is the **amount of connections between the nodes** in layer. This affects networks ability to solve given (or possibly its inability to do so) while using too few interconnections. In opposite case the network tends to produce noise in the system. So we are looking for optimal solution with *fewest possible interconnections*.

■ 2.4.2 Training

Training, or learning is unmistakably integral part of neural networks, that has to be considered, while comparing with other systems. Also, because the training phase is usually the one which requires the most time, this is a good area to make comparisons according to [38].

First criterion that comes to mind could be the **time efficiency** - how fast can the NN be trained, converging to the optimal solution. If we would discuss the fact, that learning is iterative process, and we are feeding the training data to the network in cycles, converging to expected solution; that aim should be to do so in *least epoch/iterations as possible*.

Further we can also talk about the actual real **clock time**, that would be consumed by training on given hardware. With growing complexity and depth of NNs, the required hardware capabilities are to be consider. All of this is in the end defined by the time spent, which we expect to be *as short as possible*.

Training also involves the **requirements for data set**, that would be descriptive enough for our task. This means, that we need to feed the network data of some variety, in cases of general tasks, e.g. while training NN for image compression, feeding the network only portraits of people, we couldn't expect good results while compressing (and mainly reconstructing) the image of scenery in nature. The features extracted might not be general enough for both, therefore model would have more specific application. Summing up the idea, our criterion would be *smallest possible data set*, that would fit the expected results.

Nelson and Rogers [38] describe a **computational complexity** of operations during one epoch of learning algorithm, providing example of so called Hammering and Stretching Algorithm. This method uses just one epoch to train, however requires the construction and inversion of a very large matrix, which might be detrimental during training on bigger data sets. This implies our desire of *least possible computational complexity* - the amount of operations needed.

2.4.3 Performance

The performance is logical context to be discussed as our main aim while designing learning-based solutions to our task should be based on efficiency. We already discussed time efficiency, complexity and now quality of final result, produced by our network. **How accurate** is the representation, in our case of this paper: the reconstructed image once decompressed. This might be different based on the input data and their relation to the training data set. What we mean by this, is the possibility that we are might try to process (compress) content similar to that provided during training, or vice versa; this might (or might not!) impact the resulting quality, e.g. in means of MSE or MS-SSIM in our case. Which naturally leads us to another criterion, which is the ability of network to **generalize** from the training data set. As a result, some researches include not only training data set, but also evaluation data set for purpose of accurate quality assessment (which in our case is represented by Kodak image set [26]).

We can also describe the **dependency on initial setup of weights** in our NN. Since a lot of networks start training phase with random weights, it leads to the different starting point of finding the optimal solution to our task, therefore best possible network. This implies that there's possibility of training getting stuck on a local minima, instead of being able to find the most optimal global minimum, and most effective network.

Dependency on data set order fed to NN it also to be considered. There are basically two ways to train a neural network. The first is to update the weights after calculating the error for each exemplar. The second is to accumulate the weight updates as each exemplar in the epoch is presented. At the end of the epoch (presentation of each training exemplar one time), the weight updates are made. For epoch training, the order of the presentation of exemplars will make no difference. Exemplar training can result in different results, depending on the order of the exemplar presentation. [38] Optima case would therefore consist of complete independency on order, in which the data are fed to network during training.

Specifically for image processing networks, we might focus on the final **quality range** of our system. In some cases, the trained model might not be able to achieve best possible quality in relation to target compression rate for example. Let's describe it on specific example - Agustsson et al. (2018) [1] presented GAN compression system, with main focus on low quality images, i.e. giving best possible quality outcome at extremely high compression rate. This implies that you might not get satisfying results for higher quality

compression in comparison to other systems (be it ML-based, or conventional). [14]

Somewhat similar to quality range, we might discuss **flexibility of bitrate control**. Our network might be designed in a way, that doesn't allow variability and works with fixed bitrate model(s). On the contrary [55][54] both presents *single* model, which can specify target bitrate by so called lambda parameter. Once again situational criterion, that might have a trade-off of being flexible at the cost of achievable quality. [14]

Lastly **coding unit size** gives defines the accessibility to spatial context. With bigger size of coding unit, NN give us better results; on the other hand, we must consider the representation of input data (image resolution) - since architecture is usually tailored for a specific coding unit size, which means that images that have a resolution that is not multiple of this size need to be adapted.

Chapter 3

Practical part

Practical part should focus on comparison of implementations of various ANN systems for image compression. How they score in quality amongst each other and against traditional compression methods such as JPEG, JPEG2000, BPG. Further we would like to summarize proposed classification, or rather set of criteria to be considered for classification of learning based systems.

3.0.1 Models

Twelve learning based solutions developed in recent years were chosen, available online through Tensorflow and PyTorch platforms, to be qualitatively compared with three traditional compression methods, namely **JPEG**, **JPEG2000** and **BPG**.

As for the traditional models of JPEG and JPEG2000; MATLAB implementations, through `imwrite()` function as follows:

```
% JPEG
imwrite(<input>, <output>, 'jpeg', 'quality', <qp>);
% JPEG2000
imwrite(<input>, <output>, 'jp2', 'CompressionRatio', <ratio>);
```

Free binary distribution for Windows of BPG implements codec function, that can be called from console (command prompt):

```
bpgenc -q <qp> -o <output> <input>
```

Lastly the implementation of learning based models was done through Python, specifically through platforms Tensorflow (developed by Google) and PyTorch. Following table list the encoding command and references to each model used.

| Codec | Command | Ref |
|----------------|---|------|
| HierMSE | <code>python tfci compress mbt2018-mean-mse-qp input output</code> | [35] |
| HierhMS-SSIM | <code>python tfci compress mbt2018-mean-msssim-qp input output</code> | [35] |
| HyperMSE | <code>python tfci compress bmshj2018-hyperprior-mse-qp input output</code> | [4] |
| HyperMS-SSIM | <code>python tfci compress bmshj2018-hyperprior-msssim-qp input output</code> | [4] |
| FactorMSE | <code>python tfci compress bmshj2018-factorized-mse-qp input output</code> | [4] |
| FactorMS-SSIM | <code>python tfci compress bmshj2018-factorized-msssim-qp input output</code> | [4] |
| NonlinLeaky128 | <code>python tfci compress b2018-leaky-relu-128-qp input output</code> | [3] |
| NonlinLeaky192 | <code>python tfci compress b2018-leaky-relu-192-qp input output</code> | [3] |
| NonlinGdn128 | <code>python tfci compress b2018-gdn-128-qp input output</code> | [3] |
| NonlinGdn192 | <code>python tfci compress b2018-gdn-192-qp input output</code> | [3] |
| Hific | <code>python tfci compress hific-qp input output</code> | [33] |
| ChengAnchor | <code>python3 codec.py encode cheng2020-anchor-qp -o output input</code> | [5] |

3.0.2 Metrics

We already described chosen metrics for the evaluation in section 2.3.3, so here we'll briefly describe their implementation. Most of the metrics were implemented in Matlab environment, either as part of image processing toolbox (PSNR, MS-SSIM) or freely distributed m-file scripts (PSNR-HVS-M, FSIM). Lastly the VIFP was implemented as part of Python package called Sewar. Further below is part of my evaluation code as an example of use.

```
% Evaluate the quality with full ref. objective metrics
% Feature Similarity Index (RGB)
FSIM(i-2,j-2) = FeatureSIM(ref_im, rec_im);

% Multi-scale structural similarity index measure
MSSSIM(i-2,j-2) = multissim(Y_rec,Y_ref);

% Peak signal to noise ratio (Y)
PSNR(i-2,j-2) = psnr(Y_rec,Y_ref);

% PSNR taking into account Contrast Sensitivity Function (CSF)
% and between-coefficient contrast masking of DCT basis functions.
PSNR_HVS_M(i-2,j-2) = psnrhvsm(Y_rec, Y_ref, 8);
```

As was mentioned previously, not every metric implementation works with 3-channel (RGB) images. Therefore we need to extract luminance matrix Y , which is easily done by following recommendation of ITU-R BT.709:

```
% My Matlab implementation
function Y = getLuminance(image)
R = image(:,:,1);
G = image(:,:,2);
B = image(:,:,3);
Y = R*.2126 + G*.7152 + B*.0722; % as defined in Rec. 709
```

3.0.3 Performance evaluation

Now that we have raw data from objective quality assessment, we should evaluate it and compare to subjective quality of reconstructed images available to us. For this purpose we will use specific images from the Kodak data set, representing common features.

For the start let's take **kodim19** (fig. 3.4), specifically the part with the fence (fig 3.2), representing change in spatial frequency. This is common feature problematic during compression, thanks to the aliasing. Therefore could be source of compression artifacts.



Figure 3.1: Visual comparison of reconstructed image Kodim19 from Kodak, $\text{bpp} \approx 0.11$

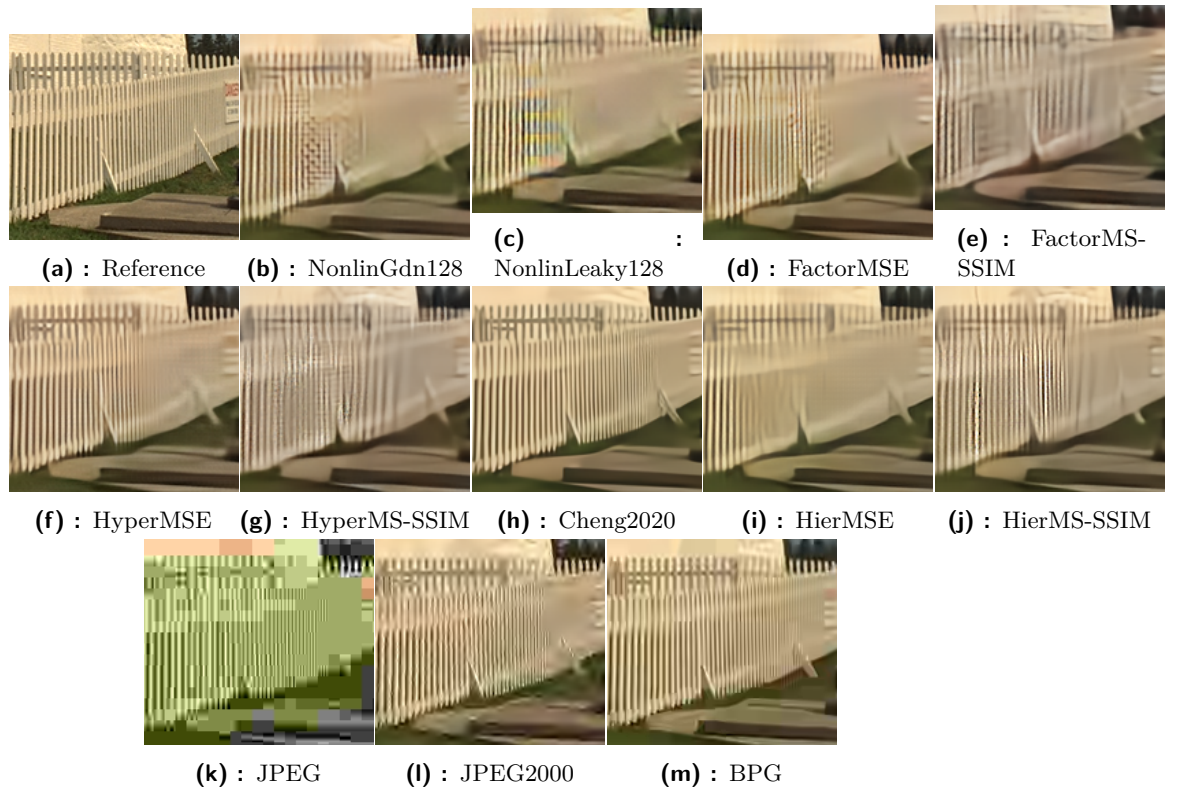


Figure 3.2: Area of interest (fence), cropped from reconstructed image Kodim19 from Kodak, $\text{bpp} \approx 0.11$

By purely subjective assessment, we could see that models **HyperMSE** gives somewhat satisfying results, although **BPG** and **Cheng2020** definitely prevail, showing next to none aliasing related artifacts. Other models either lose sharpness fairly fast, or creates specific repeated pattern in the area. Model **NonlinLeaky128** even creates moiré, specific colour aberration common with aliasing. Legacy **JPEG** proves to be almost unusable on such low quality settings (0.11 bpp), thanks to heavy loss in colour resolution, and blocky structure. This will prove to be common for JPEG compression with high compression ratio setting.

Further we shall take a look at **Kodim20**, and crop out of logo from the reconstructed images. Here we can explore the readability of signs, logos etc. Once again, subjectively we can define reconstructed image from model **HyperMSE**, although specific artifact - repeated blocky pattern on whole image (as seen on 3.3. **JPEG2000** also satisfies the criteria of readability, although compression artifacts common for JPEG2000 hinders the overall quality on high compression levels. Loss of contrast on areas with high brightness (the sky) is noticeable using **FactorMS-SSIM** and **HyperMS-SSIM** as well. Once again **Cheng2020** proves to be most satisfying subjective quality.

Let us compare how this subjectively assessed performance correlates with objectively measure quality of chosen metrics. The important metrics this



Figure 3.3: Visualization of block pattern artifact from HyperMSE model.

work focuses on are VIFP, FSIM and PSNR-HVS-M, since they bring superior results correlating with subjective assessments while using learning based methods for compression.

Going back to our image **Kodim19** (lighthouse) we have measured objective quality score as seen on 3.6 and 3.7. For metrics PSNR-HVS-M and VIFP, we can clearly see that model Cheng2020 and BPG are superior to other models. This supports our conclusion from subjective assessment. On the opposite end, we can see that JPEG doesn't its performance is sub par, be it in contrast to learning based metrics or other traditional methods (which isn't any surprise).

We can further derive that most "dual" models that have mse-based and ms-ssim-based variants give very similar results. Lastly JPEG2000 scores average among all mentioned models.

Although **FSIM** scores suggest, that two ms-ssim-based models (namely **HierMS-SSIM** and **HyperMS-SSIM**) achieve similar or even superior quality to Cheng2020 and BPG. This might mean that FSIM scores doesn't take into account specific artifacts mentioned before; in contrast to subjective assessment, this could be contested.

Lastly MS-SSIM heavily scores well on JPEG for bpp over 0,2. That is highly in contradiction with subjective assessment, since JPEG images around this bitrate are still quite blocky which highly influences the quality.

Now let's compare the objective scores given for Kodim20, therefore see how models and metrics relate contextually also. Looking at the 3.8 and 3.9, we can see some similar trends as well as differences in objective scores. Firstly, models **Cheng2020** and **BPG** proves to be superior in all metrics

used. On the other hand the **PSNR-HVS-M** metric shows different trend for all ms-ssim-based models. This might be related to big homogeneous part of the image with lack of details (sky), and correlates to our assessment of loss of contrast for these particular models. This is further supported by simple PSNR metric.

Once again, the MS-SSIM favours JPEG for bpp over 0,2. This might mean, that MS-SSIM might not be completely viable when assessing quality of learning based systems.

Rest of the models follows trends that we stated before, and we don't see much differences throughout all metrics.



Figure 3.4: Visual comparison of reconstructed image Kodim20 from Kodak, $\text{bpp} \approx 0.11$

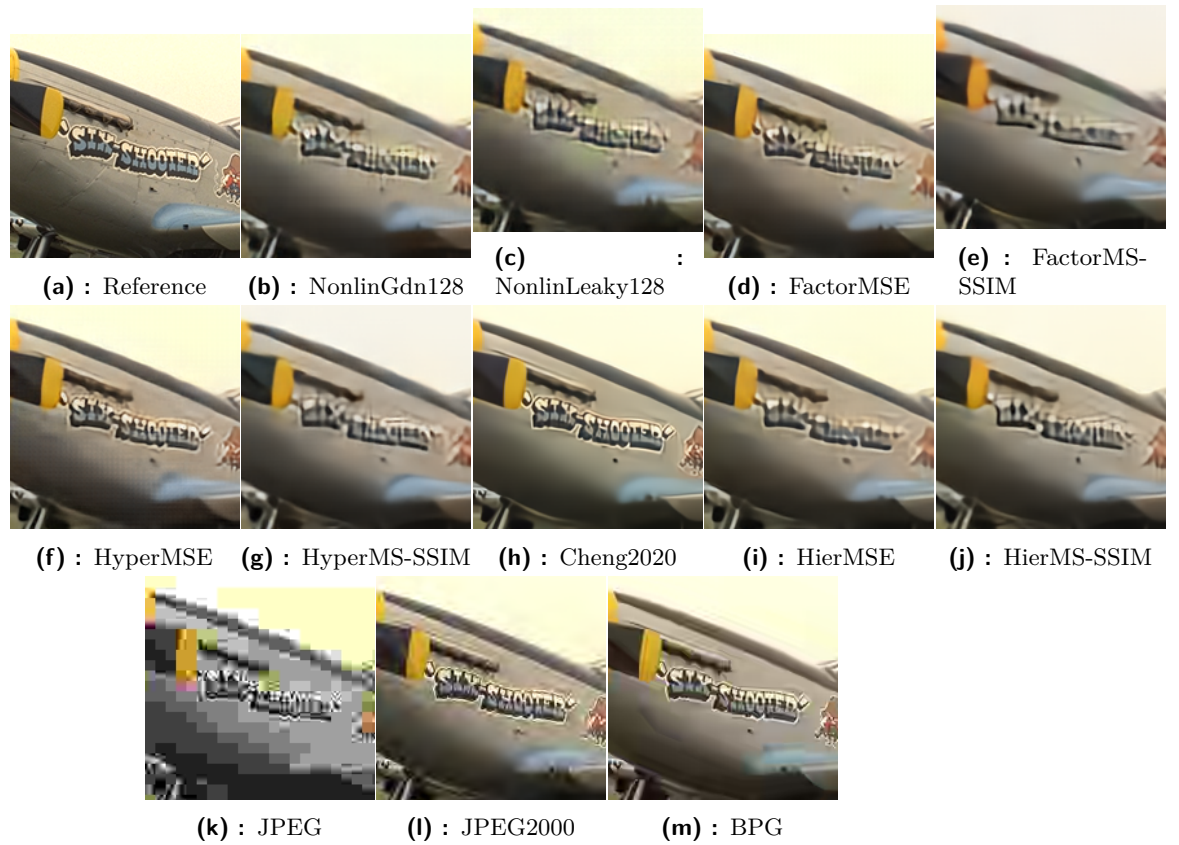
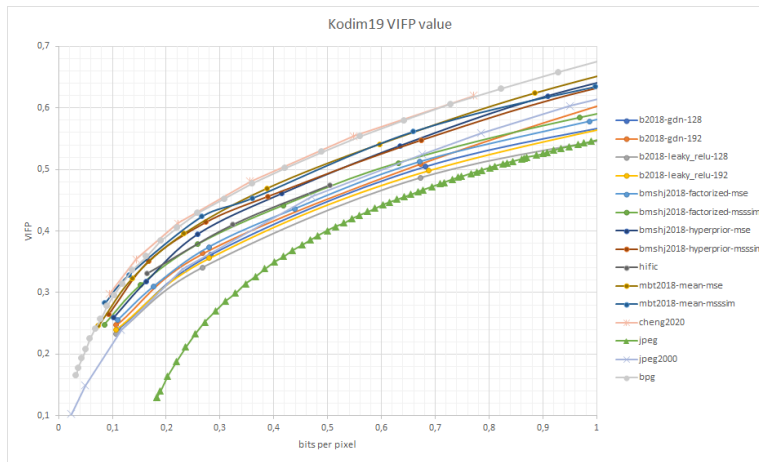
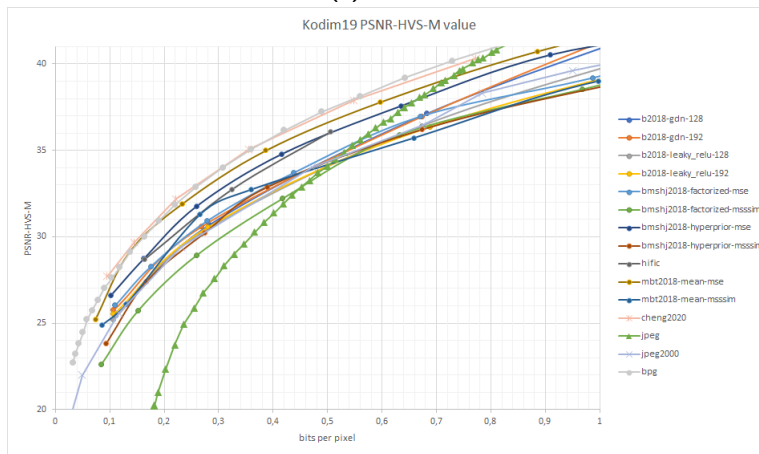


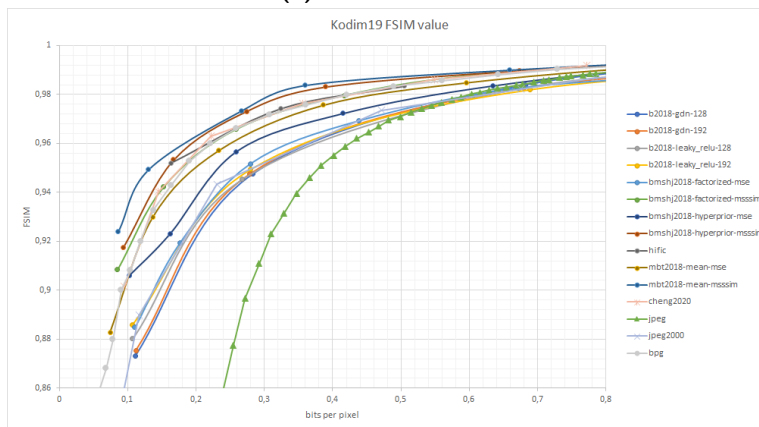
Figure 3.5: Area of interest (logo), cropped from reconstructed image Kodim20 from Kodak, $\text{bpp} \approx 0.11$



(a) : VIFP

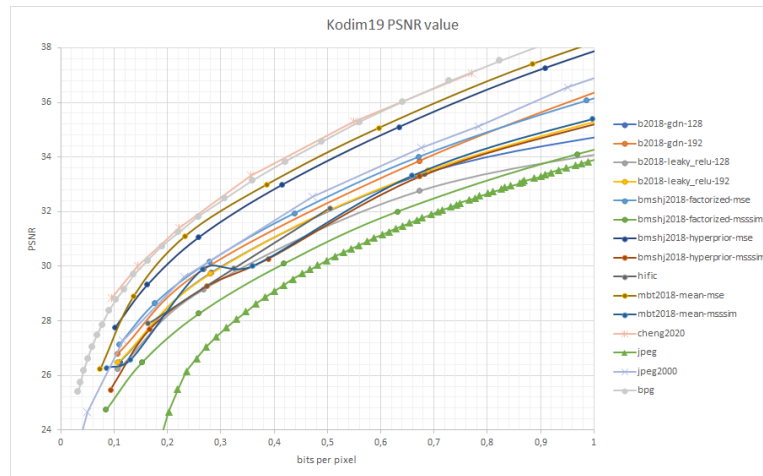


(b) : PSNR-HVS-M

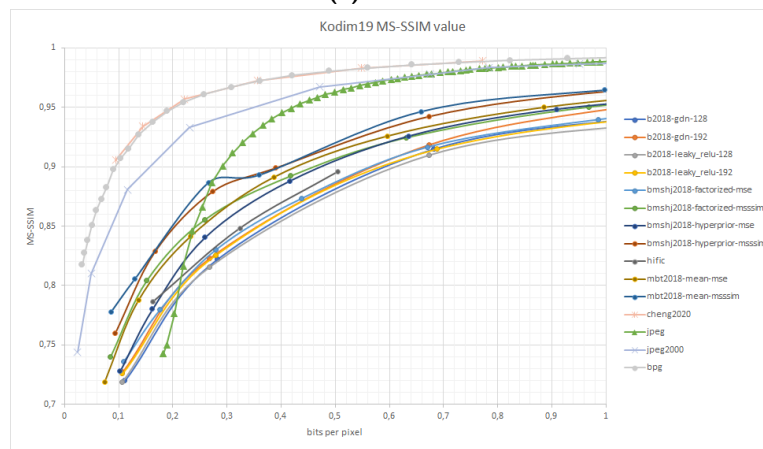


(c) : FSIM

Figure 3.6: Measured scores for various objective quality metrics (Kodim19).

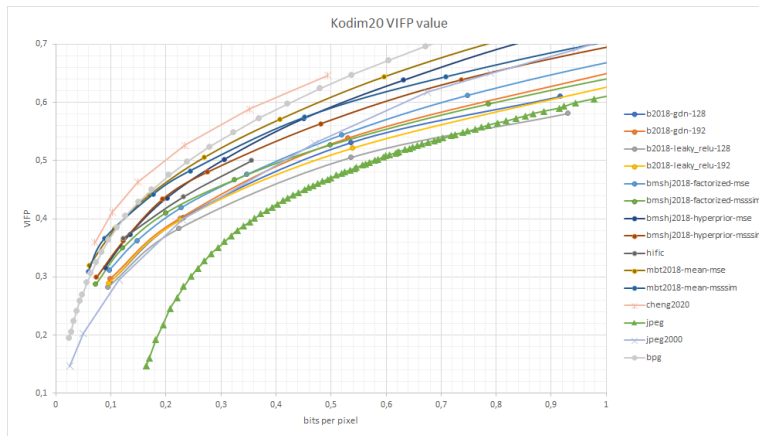


(a) : PSNR

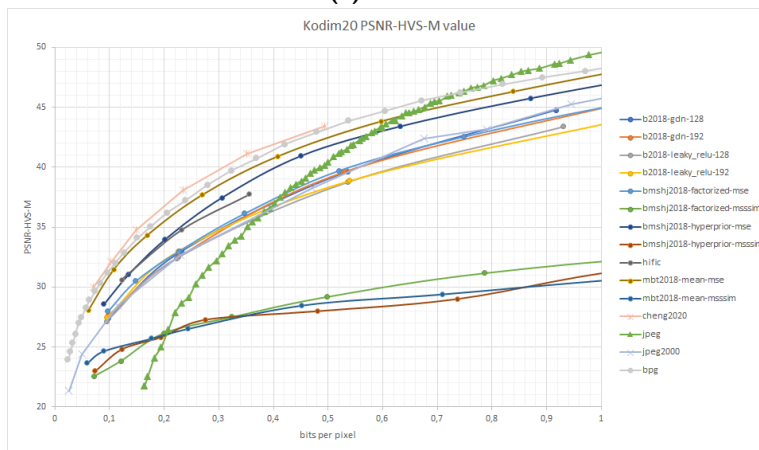


(b) : MS-SSIM

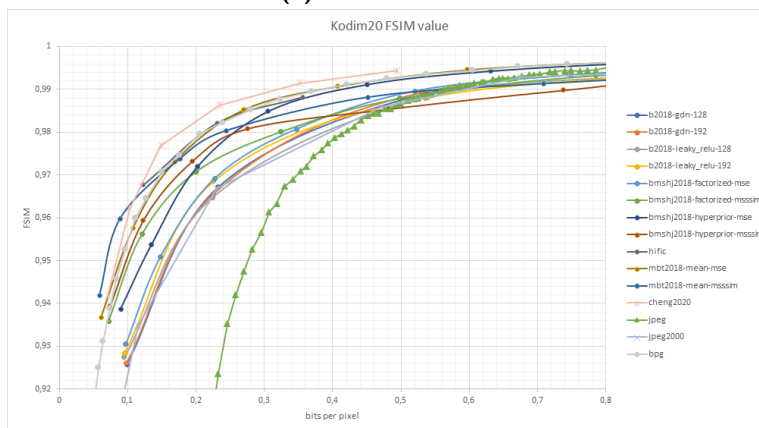
Figure 3.7: Measured scores for various objective quality metrics (Kodim19).



(a) : VIFP

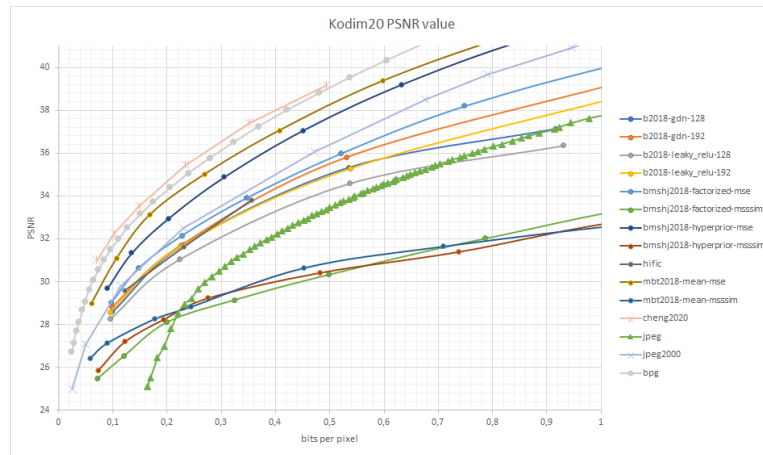


(b) : PSNR-HVS-M

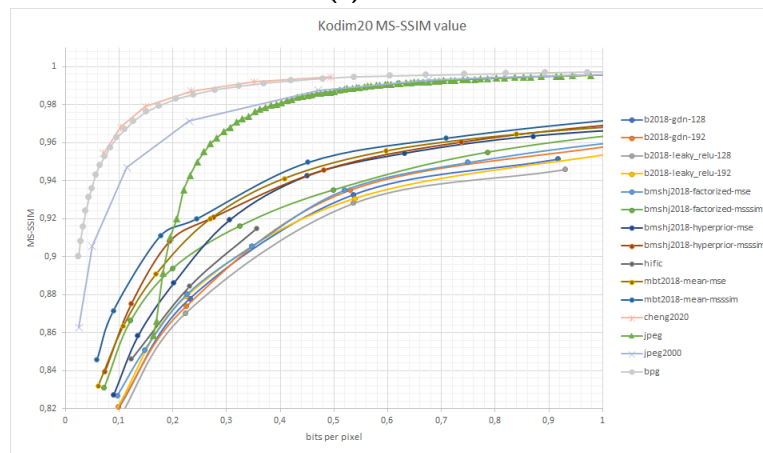


(c) : FSIM

Figure 3.8: Measured scores for various objective quality metrics (Kodim20).



(a) : PSNR



(b) : MS-SSIM

Figure 3.9: Measured scores for various objective quality metrics (Kodim20).

Chapter 4

Discussion and conclusions

The results of our research clearly show, that ML based compression methods proves to be step forward in efficient compression of visual data. It is very promising field, with ability to already qualitatively exceed the conventional methods of image processing. However alongside with research in new compression methods, there is also immense need for state-of-the-art objective quality metrics. As proven by the presented results of various metrics used, the modified (or completely new) metrics for the purpose of learning based methods are needed, to give us precise score correlated to the subjective scores. FSIM, VIFP as well as PSNR-HVS-M proved to be right choice for such assessment of performance.

As for the output quality of learning based models: quality of the reconstructed images are on par, in some cases superior to the traditional codecs. This is especially true for low bitrates, high compression setups. Model Cheng2020 implemented through CompressAI (PyTorch library) shows very promising results among the others, which is to no surprise as it's the most recent method used in this work.

This is strictly speaking from performance point of view. The conventional methods are also characterized by their easy of use, after the years of research and development. They are flexible as to the format of the input, quality range provided and implementations in most of the modern systems. Same can't be told for ML-based systems; as for now at least. Another thing to consider is possibility of distribution, and hardware requirements on the consumer/user side. Last but not least, the learning/training of new models is trivial matter still, and appears to be quite time-consuming, without proper high-end hardware (GPU). Considering all of these milestones, I'm curious to see how future of standardization by efforts of JPEG-AI will be.

As for the classification (taxonomy) of ML-based system, there's plethora of criteria to look for, extensively described in 2.4. Also many of these criteria are connected and interlaced, sometimes forcing us to make trade-offs, e.g. efficiency vs. flexibility etc. Also not every criterion is applicable to all problem areas; this implies a lot of decision and criteria application has to be done in advance with fore knowledge.

This paper tries to navigate the ever-growing field ANNs and their application in image compression; that is specifically in comparison to conventional

compression methods. Practical part there is focuses on this fact is holding possibilities for further improvement, and would be worth of more extensive research. Specifically considering the platforms used for NN implementations - this paper solely focused on methods based on Python language and on TensorFlow platform and PyTorch. Other (still Python-based) platform are available (such as) and would be worth of deeper research, such as Theano, OpenCV or Keras; they might provide valuable data to be compared with those provided in this paper.



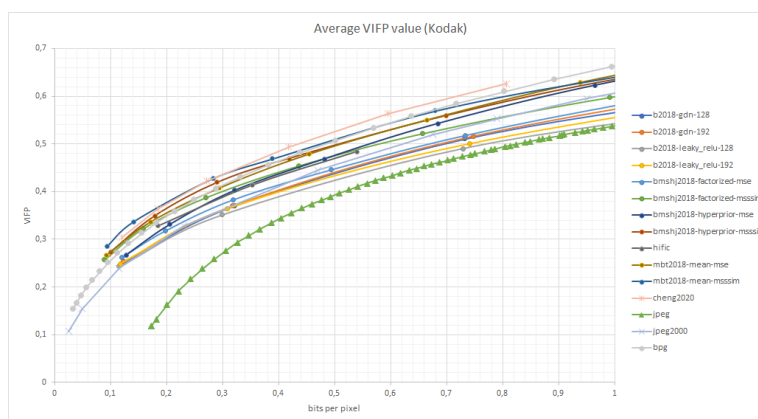
Appendix A

List of electronic attachments

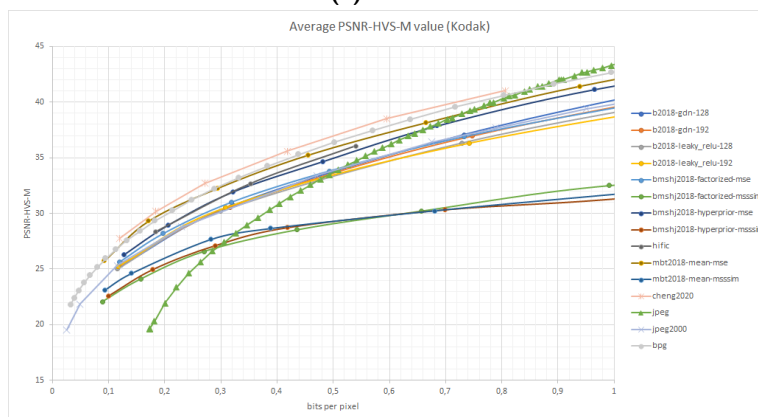
- `evaluation.m`: evaluation MATLAB script for JPEG and JPEG2000 compression, with resulting MS-SSIM, PSNR, FSIM, PSNR-HVS-M scores, as well as BPP of compressed files
- `codec.m`: MATLAB script for compression and reconstruction of image data, with JPEG and JPEG2000
- `evaluation.py`: Python script for VIFP score evaluation
- `codec.py`: Python script for calling various compression models accessible through `tfc`.py
- `tfc`.py: python script implementing various compression methods, downloading related models (availability changes during time); by Ballé et al.
- `resultsFull.xlsx`: Excel file containing all evaluation scores
- `getLuminance.m`: MATLAB function to get Y component from RGB image

Appendix B

Average objective scores for Kodak dataset



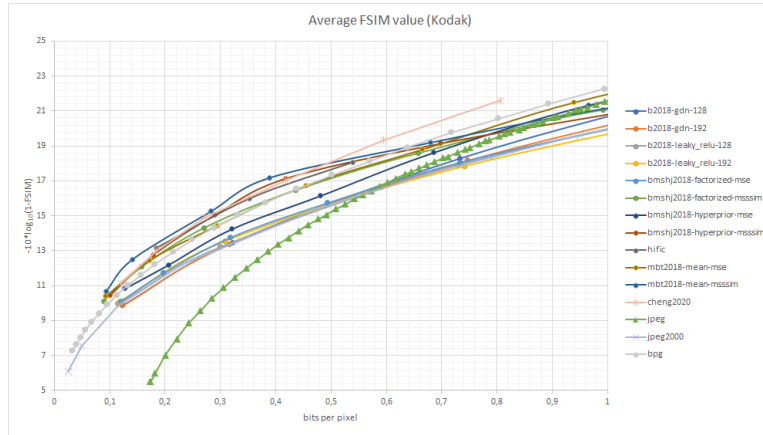
(a) : VIFP



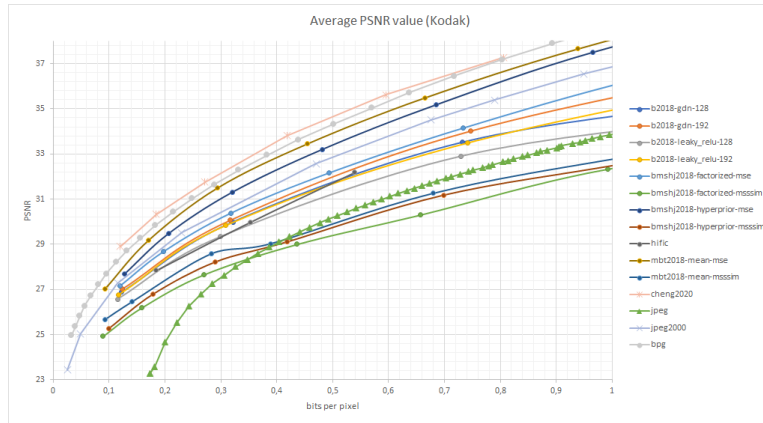
(b) : PSNR-HVS-M

Figure B.1: Average objective scores over Kodak data set.

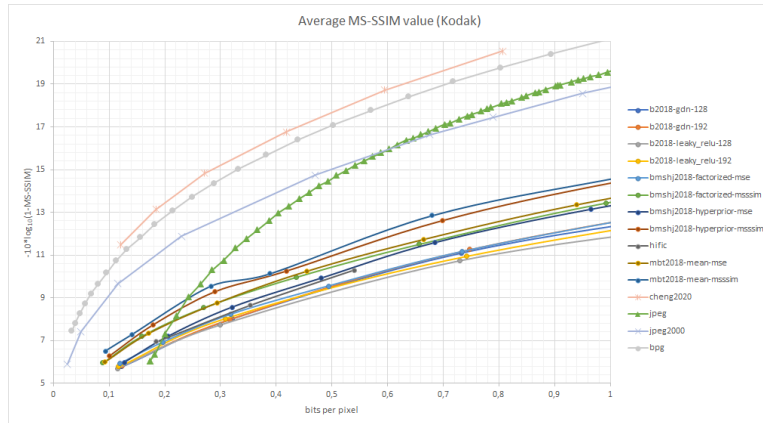
B. Average objective scores for Kodak dataset



(a) : FSIM



(b) : PSNR



(c) : MS-SSIM

Figure B.2: Average objective scores over Kodak data set.



Bibliography

- [1] Eirikur Agustsson et al. *Generative Adversarial Networks for Extreme Learned Image Compression*. 2019. arXiv: 1804.02958 [cs.CV].
- [2] N. Asuni and A. Giachetti. “TESTIMAGES: A large-scale archive for testing visual devices and basic image processing algorithms (SAMPLING 1200 RGB set)”. In: *STAG: Smart Tools and Apps for Graphics*. 2014. URL: https://sourceforge.net/projects/testimages/files/OLD/OLD_SAMPLING/testimages.zip.
- [3] Johannes Ballé. “Efficient nonlinear transforms for lossy image compression”. In: *2018 Picture Coding Symposium (PCS)*. IEEE. 2018, pp. 248–252.
- [4] Johannes Ballé et al. *Variational image compression with a scale hyperprior*. 2018. arXiv: 1802.01436 [eess.IV].
- [5] Jean Bégaint et al. *CompressAI: a PyTorch library and evaluation platform for end-to-end compression research*. 2020. arXiv: 2011.03029 [cs.CV].
- [6] H. Bourlard and Y. Kamp. “Auto-association by multilayer perceptrons and singular value decomposition”. In: *Biological Cybernetics* 59 (2004), pp. 291–294.
- [7] Leo Breiman. “Random forests”. In: *Machine learning* 45.1 (2001), pp. 5–32.
- [8] Yung-Yao Chen et al. “Design and Implementation of Cloud Analytics-Assisted Smart Power Meters Considering Advanced Artificial Intelligence as Edge Analytics in Demand-Side Management for Smart Homes”. In: *Sensors* 19 (May 2019), p. 2047. DOI: 10.3390/s19092047.
- [9] *Common ML Problems / Introduction to Machine Learning Problem Framing*. URL: <https://developers.google.com/machine-learning/problem-framing/cases>.
- [10] Corinna Cortes and Vladimir Vapnik. “Support-vector networks”. In: *Machine learning* 20.3 (1995), pp. 273–297.
- [11] Thomas Cover and Peter Hart. “Nearest neighbor pattern classification”. In: *IEEE transactions on information theory* 13.1 (1967), pp. 21–27.

- [28] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. “Deep learning”. In: *nature* 521.7553 (2015), pp. 436–444.
- [29] Yann Lecun et al. “Gradient-Based Learning Applied to Document Recognition”. In: *Proceedings of the IEEE* 86 (Dec. 1998), pp. 2278–2324. DOI: 10.1109/5.726791.
- [30] Dong Liu et al. “Deep learning-based video coding: A review and a case study”. In: *ACM Computing Surveys (CSUR)* 53.1 (2020), pp. 1–35.
- [31] Siwei Ma et al. “Image and video compression with neural networks: A review”. In: *IEEE Transactions on Circuits and Systems for Video Technology* 30.6 (2019), pp. 1683–1698.
- [32] *Machine Learning Glossary | Google Developers*. URL: <https://developers.google.com/machine-learning/glossary>.
- [33] Fabian Mentzer et al. “High-fidelity generative image compression”. In: *arXiv preprint arXiv:2006.09965* (2020).
- [34] Donald Michie. ““Memo” functions and machine learning”. In: *Nature* 218.5136 (1968), pp. 19–22.
- [35] David Minnen, Johannes Ballé, and George Toderici. *Joint Autoregressive and Hierarchical Priors for Learned Image Compression*. 2018. arXiv: 1809.02736 [cs.CV].
- [36] Tom M. Mitchell. *Machine Learning*. New York: McGraw-Hill, 1997. ISBN: 978-0-07-042807-2.
- [37] Radford M. Neal and Peter Dayan. “Factor Analysis Using Delta-Rule Wake-Sleep Learning”. In: *Neural Comput.* 9.8 (Nov. 1997), pp. 1781–1803. ISSN: 0899-7667. DOI: 10.1162/neco.1997.9.8.1781. URL: <https://doi.org/10.1162/neco.1997.9.8.1781>.
- [38] D. E. Nelson and S. K. Rogers. “A taxonomy of neural network optimality”. In: *Proceedings of the IEEE 1992 National Aerospace and Electronics Conference@m_NAECON 1992*. 1992, 894–899 vol.3.
- [39] M.A. Nielsen. *Neural Networks and Deep Learning*. Determination Press, 2015. URL: <http://neuralnetworksanddeeplearning.com/>.
- [40] Jo Plumridge. *Are You Losing Too Much of Your Photograph to Compression?* July 2020. URL: <https://www.lifewire.com/the-effect-of-compression-on-photographs-493726>.
- [41] Nikolay Ponomarenko et al. “On between-coefficient contrast masking of DCT basis functions”. In: *Proc of the 3rd Int Workshop on Video Processing and Quality Metrics for Consumer Electronics* (Jan. 2007).
- [42] David Poole, Alan Mackworth, and Randy Goebel. “Computational Intelligence”. In: (1998).
- [43] Bin Qian et al. *Orchestrating the Development Lifecycle of Machine Learning-Based IoT Applications: A Taxonomy and Survey*. 2019. arXiv: 1910.05433 [cs.DC].

