

Foundations and Trends® in Computer Graphics
and Vision

An Introduction to Neural Data Compression

Suggested Citation: Yibo Yang, Stephan Mandt and Lucas Theis (2023), “An Introduction to Neural Data Compression”, Foundations and Trends® in Computer Graphics and Vision: Vol. 15, No. 2, pp 113–200. DOI: 10.1561/0600000107.

Yibo Yang

University of California, Irvine
yibo.yang@uci.edu

Stephan Mandt

University of California, Irvine
mandt@uci.edu

Lucas Theis

Google Research
theis@google.com

This article may be used only for the purpose of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval.

now

the essence of knowledge

Boston — Delft

Contents

1	Introduction	114
2	Lossless Compression	119
2.1	Background	119
2.2	Autoregressive models	126
2.3	Latent variable models	128
2.4	Invertible flows and other models	134
3	Lossy Compression	135
3.1	Background	136
3.2	Neural lossy compression	143
3.3	Learned quantization and rate control	148
3.4	Compression without quantization	158
3.5	Perceptual losses	163
3.6	Task-oriented compression	171
3.7	Video compression	172
4	Discussion and Open Problems	175
	Acknowledgements	178
	References	179

An Introduction to Neural Data Compression

Yibo Yang¹, Stephan Mandt² and Lucas Theis³

¹*University of California, Irvine, USA; yibo.yang@uci.edu*

²*University of California, Irvine, USA; mandt@uci.edu*

³*Google Research, USA; theis@google.com*

ABSTRACT

Neural compression is the application of neural networks and other machine learning methods to data compression. Recent advances in statistical machine learning have opened up new possibilities for data compression, allowing compression algorithms to be learned end-to-end from data using powerful generative models such as normalizing flows, variational autoencoders, diffusion probabilistic models, and generative adversarial networks. This monograph aims to introduce this field of research to a broader machine learning audience by reviewing the necessary background in information theory (e.g., entropy coding, rate-distortion theory) and computer vision (e.g., image quality assessment, perceptual metrics), and providing a curated guide through the essential ideas and methods in the literature thus far.

Yibo Yang, Stephan Mandt and Lucas Theis (2023), “An Introduction to Neural Data Compression”, Foundations and Trends® in Computer Graphics and Vision: Vol. 15, No. 2, pp 113–200. DOI: 10.1561/0600000107.

©2023 Y. Yang, S. Mandt and L. Theis

1

Introduction

The goal of data compression is to reduce the number of bits needed to represent useful information. *Neural*, or *learned* compression, is the application of neural networks and related machine learning techniques to this task. This monograph aims to serve as an entry point for machine learning researchers interested in compression by reviewing the prerequisite background and representative methods in neural compression.

The basic idea of learning-based data compression has long existed in various forms before the current era of deep learning [224][154][37][60]. Many of the tools and techniques for neural compression, especially for images, also draw on a rich history of learning-based approaches in computer vision. Indeed, many problems in image processing and restoration can be viewed as lossy image compression; e.g., image super-resolution can be solved by learning a decoder for a fixed encoder (the image downsampling process) [49][105]. In fact, neural networks have already been applied to image compression in the late 1980s and 1990s [170][61], and even an early review article [96] exists. Compared to early work, modern methods differ markedly in their scale, neural architectures, and encoding schemes.



Figure 1.1: Compression as generative modeling. *Left:* A sample drawn from the probabilistic model underlying JPEG, which betrays an assumption of independence among neighboring 8 by 8 pixel blocks (except for the DC components within each row). *Right:* A sample generated by a recent neural compression model by Minnen *et al.* [132].

Current research in neural compression is heavily inspired by advances in deep generative modeling, such as GANs [65], VAEs [99][151], normalizing flows [104], and autoregressive models [180][140]. While these models allow us to capture complex data distributions from samples (a key to neural compression), the research tends to focus on generating realistic data samples [139][142] or achieving high data log-density [151][101], objectives not necessarily aligned with data compression.

Arguably the first work exploring deep generative models for data compression appeared in 2016 [70], and the topic of neural compression has grown considerably since then. Multiple researchers have identified connections between variational inference and lossless [59][118] as well as lossy [12][184][6][209] compression. This monograph hopes to further facilitate such exchange between these fields, raising awareness of compression as a fruitful application of generative modeling along with the associated challenges.

Instead of surveying the vast literature, we aim to cover the essential concepts and methods in neural compression, with a reader in mind who is versed in machine learning but not necessarily data compression. We hope to complement existing surveys that have a more specialized or applied focus [10][117][111] by highlighting the connections to generative modeling and machine learning in general. In most of this monograph, we make essentially no assumption on the data other than that it

is independently and identically distributed (i.i.d.), a typical setting for machine learning and statistics. We center our discussions around image compression, where most neural compression methods were first developed, but the basic ideas we present here are data agnostic. Towards the end, in Section 3.7, we lift the i.i.d. assumption and consider video compression, which can be seen as an extension of the existing ideas along the temporal dimension.

Neural compression can ease the development and optimization of data compression algorithms in a data-driven fashion. This can be especially useful for new or domain-specific data types, such as VR/AR content or scientific data, where developing custom codecs may otherwise be expensive. Indeed, learning-based approaches are being applied to emerging data types, such as point clouds [147][72][89], implicit 3D surfaces [178], and neural radiance fields [22]. Effectively compressing such data may require new neural architectures [178] and/or domain knowledge to convert the data into neural-network-friendly representations [89]. However, the essential ideas and techniques introduced here for reducing the entropy, or bit-rate cost, of learned representations remain the same.

JPEG [92] serves as a good motivating example of the lossy compression pipeline (depicted in Figure 1.2). First introduced in 1992, it is still one of the most widely used image compression standards [90]. At the heart of JPEG are linear mappings which losslessly transform pixels into coefficients and back. The coefficients are first quantized to integers, incurring some information loss. Then they are further compressed losslessly by a combination of run-length encoding and entropy coding (the latter is discussed in Section 2.1.1).

The linear portion of the encoding process consists of several steps. First, each pixel is transformed from RGB to YCC coefficients consisting of a luma component (Y) and two color components (C). After this color transform, each channel is treated independently, and optional downsampling is applied to the color channels. Next, each channel is divided into 8×8 pixel blocks, and each block independently undergoes a *discrete cosine transform* (DCT). The transform coefficients are then

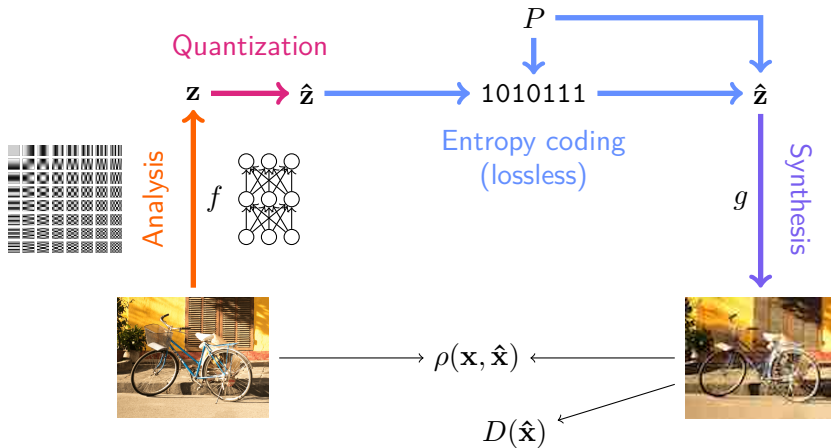


Figure 1.2: A typical pipeline in both neural and classical lossy image compression. An encoder transformation f (for example, the DCT or a neural network) maps images to coefficients \mathbf{z} , which are first quantized to $\hat{\mathbf{z}}$, and then entropy encoded into bits using an entropy model P . A reconstruction $\hat{\mathbf{x}}$ is obtained using a decoder g that aims for a small distortion ρ between the data \mathbf{x} and its lossy reconstruction $\hat{\mathbf{x}}$. In addition, neural compression can also involve an adversarial critic D , encouraging realism and high perceptual quality.

linearly scaled and finally rounded to integers. Given an image \mathbf{x} , the encoder thus performs

$$\hat{\mathbf{z}} = \lfloor \mathbf{D}\mathbf{A}\mathbf{C}\mathbf{x} \rfloor, \quad (1.1)$$

where \mathbf{C} is the pixelwise color transform, \mathbf{A} is the block- and channelwise DCT, and \mathbf{D} is a diagonal matrix scaling the coefficients. The decoder applies the transforms in reverse,

$$\hat{\mathbf{x}} = \mathbf{C}^{-1}\mathbf{A}^{\top}\mathbf{D}^{-1}\hat{\mathbf{z}}. \quad (1.2)$$

Readers familiar with machine learning will be reminded of autoencoders [29][158] and it is natural to consider learned neural networks in place of the linear transforms. As we will see later, there are indeed close connections between lossy compression and variational autoencoders (VAEs) [12][184][6][211], though other generative models have a role to play as well. What we call “coefficients” in the context of compression are often called “latent variables” in the context of generative models.

Like generative models, JPEG defines a probability distribution over coefficients which represents assumptions about the latent representation. Just as in VAEs, we can use this distribution to draw samples from the model underlying JPEG, with an example shown in Figure 1.1.

Overview. This introduction is organized into two main parts, lossless (Section 2) and lossy (Section 3) compression, with the latter relying on the former for compressing lossy representations of the data (see Figure 1.2). We begin by reviewing basic *coding theory* (Section 2.1), and learn how we can turn the problem of lossless compression into learning a discrete data distribution, with the help of *entropy-coding*. For this to work in practice, we decompose the potentially high-dimensional data distribution using tools from generative modeling, including *autoregressive models* (Section 2.2), *latent-variable models*, (Section 2.3), and other models (Section 2.4). Each model class differs in its compatibility with different entropy-coding algorithms, and offers a different trade-off between the compression bit-rate and computational efficiency. *Lossy* compression introduces additional desiderata, the most common being the *distortion* of reconstructions, based on which the classical *rate-distortion theory* and algorithms such as *vector quantization* and *transform coding* are reviewed (Section 3.1). We then introduce *neural lossy compression* as a natural extension of transform coding (Section 3.2) and discuss the techniques necessary for end-to-end learning of quantized representations (Section 3.3), as well as lossy compression schemes that attempt to bypass quantization (Section 3.4). We then explore additional desiderata, such as the *perceptual quality* of reconstructions (Section 3.5), and the usefulness of learned representations for downstream tasks (Section 3.6), before briefly reviewing video compression (Section 3.7). Finally, we conclude in Section 4 with the challenges and open problems in neural compression that may drive its future advances.

2

Lossless Compression

Lossless compression aims to represent data with as few bits as possible such that the data can be reconstructed perfectly. The basic recipe is to first build a probabilistic model of the data, and then feed its probabilities into a so-called entropy coding scheme which converts data into compact bit-strings. This precludes non-likelihood-based models such as *generative adversarial networks* (GANs) [65], from which it is hard to derive probabilities and which would often assign zero probability to data even if we could evaluate them.

2.1 Background

In the first few subsections, we review the basic concepts and algorithms for *entropy coding*, which is the core interface between lossless compression and data modeling. Then in Section 2.1.6, we discuss a commonly used modeling technique for representing a discrete distribution with a density.

2.1.1 Entropy coding

We call a sequence of outcomes of a discrete random variable a *message*, and *entropy coding* is a way to achieve lossless compression of a message. The basic idea, as embodied by *symbol codes*, is to replace commonly occurring symbols with short codewords and rare symbols with longer codewords, thereby reducing the message’s overall length. Morse code implements this idea by representing the common letter “e” with a single so-called *dit* and the less frequent letter “s” using 3 *dits*. To be able to distinguish between the message “eee” and the message “s”, Morse code additionally requires pauses between encoded letters. These extra markers can be avoided by instead using a *prefix-free* code where no codeword is the prefix of any other codeword [39].

By Shannon’s source coding theorem [163], the codeword length of an optimal prefix-free code is approximately the negative logarithm of the codeword’s probability. Shannon also showed that the expected message length of an optimal prefix-free code is close to the *entropy* of the message.

2.1.2 Entropy and information

The entropy¹ of a discrete random variable $\mathbf{X} \sim P$ is a measure of uncertainty about its outcomes. It is based on the concept of *surprise*, or *information content*, which can be defined as the negative log-probability of an outcome $-\log_2 P(\mathbf{x})$. Entropy, by definition, is the expected value of surprise,

$$H[\mathbf{X}] = \mathbb{E}_{\mathbf{x} \sim P}[-\log_2 P(\mathbf{x})]. \quad (2.1)$$

For example, the entropy of a fair coin flip is 1 bit, and the entropy of a biased coin flip approaches 0 bits as the probability of one of the two outcomes approaches 1.

When we losslessly compress data with a prefix-code, the entropy is the minimum number of bits required on average. More precisely,

¹Entropy was originally defined as a thermodynamic concept by Clausius [38] and later Boltzmann [27], while the information theoretic entropy was introduced by Shannon [163].

the expected message length $|C(\mathbf{x})|$ of an optimal prefix-free code C is bounded as

$$H[\mathbf{X}] \leq \mathbb{E}[|C(\mathbf{X})|] \leq H[\mathbf{X}] + 1. \quad (2.2)$$

In practice, we typically do not know the data distribution P but need to approximate it with a distribution Q . We usually estimate Q by maximum-likelihood, or equivalently, minimizing the *cross-entropy* between P and Q , defined as,

$$H[P, Q] = \mathbb{E}_{\mathbf{x} \sim P}[-\log_2 Q(\mathbf{x})]. \quad (2.3)$$

This is justified from a compression perspective, as the cross-entropy captures the average number of bits needed to code samples from P using a code optimized for Q . The cross-entropy is always at least as large as the entropy, and the gap between the two is called *relative entropy* or Kullback-Leibler (KL) divergence, defined as

$$d_{\text{KL}}[P \parallel Q] = \mathbb{E}_{\mathbf{x} \sim P}[-\log_2 Q(\mathbf{x}) + \log_2 P(\mathbf{x})]. \quad (2.4)$$

It is always positive unless $P = Q$ and serves as an asymmetric distance measure between the two distributions.

Based on these foundational concepts, we will review a few entropy coding schemes next.

2.1.3 Huffman coding

Perhaps the most well-known prefix-free symbol coding approach is Huffman coding [91]. In a nutshell, given a distribution over symbols, the algorithm assigns unique binary codewords to every symbol by building a so-called Huffman tree (Figure 2.1). The leaf nodes of the tree are associated with the symbols. A Huffman tree can be recursively grown from the leaves to the root by successively merging the two nodes that have the lowest probabilities under the data distribution P . By traversing the tree from the root to a leaf, the sequence of branching directions (“0” for left, “1” for right) assigns a unique binary codeword to each symbol, with the more frequent symbols receiving shorter codewords. Encoding and decoding are simple lookup operations with complexity linear in the tree depth.

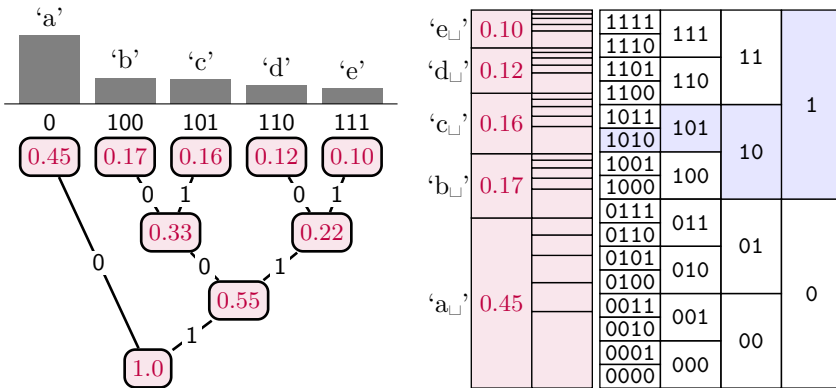


Figure 2.1: *Left:* A Huffman tree over five symbols. Purple numbers in boxes denote probabilities, while binary strings above the leaf nodes correspond to the codewords assigned by the Huffman algorithm. *Right:* Arithmetic coding of a message of length two. Each message corresponds to an interval whose length is proportional to the message’s probability. As an example, any real number inside the interval corresponding to the message ‘ca’ begins with 0.1010... when written in binary form so that we can use 1010 to uniquely encode it. The message ‘aa’ would be encoded as 000.

Huffman coding assigns a codeword with length at most $\lceil -\log_2 P(\mathbf{x}) \rceil$ for each symbol \mathbf{x} , and can be shown to be optimal among prefix codes [39]. However, as log-probabilities are generally not integers, Huffman coding (as with all so-called *symbol codes*) incurs an overhead of up to 1 bit per symbol, relative to the information content. For example, a heavily biased coin with an entropy close to 0 will still require 1 bit to encode. When compressing a sequence of symbols, *streaming codes* (see below) allow us to do so more efficiently than symbol codes by limiting the overhead to 1 bit for an entire *message* rather than each symbol.

2.1.4 Arithmetic coding

Arithmetic coding [154], also known as range coding [122], is our first example of a *streaming code*. Streaming codes differ from symbol codes in that they assign codewords to entire messages and individual symbols do not have unique codewords. Streaming codes amortize the coding cost’s overhead across the whole sequence of symbols and are therefore able to get closer to the entropy.

The basic idea of arithmetic coding is to associate each symbol with a subinterval of the interval $[0, 1]$ such that the subinterval's length equals the symbol's probability. This procedure can be generalized towards encoding symbols in sequence. For two symbols, the second symbol is no longer considered a subinterval of $[0, 1]$ but a subinterval of the previous symbol's interval such that the second subinterval's length equals the two symbols' joint probability. For example, if symbol "a" has probability $1/5$ and is associated with $[0, 0.2)$, the sequence "aa" would be encoded as $[0, 0.04)$. This procedure can be iterated for sequences of arbitrary length, leading to a sequence of subintervals of decreasing size that contain each other. Any real-valued number contained in the final subinterval allows us to uniquely reconstruct the sequence of symbols and therefore encodes the entire sequence (Figure 2.1).

In practice, one picks a representative number which can be represented with the smallest number of bits while uniquely identifying the interval. By construction, the size of the final sub-interval is the product of all previous symbols' conditional probabilities, hence the probability of the sequence. Intuitively, the interval's length (i.e., the sequence's joint probability) determines the number of relevant digits of the number representing the interval, thus the code length equals the log-probability of the sequence as required for entropy codes. If \mathbf{x}^n is a message of length n , and $C(\mathbf{x}^n)$ the codeword assigned by arithmetic coding using an assumed probability distribution Q , then $|C(\mathbf{x}^n)| < \lceil -\log_2 Q(\mathbf{x}^n) \rceil + 1$ [26]. Thus if the true distribution of messages is P , then the average code length is at most $H[P, Q] + 2$. The advantage of arithmetic coding over Huffman coding is that the excess bits needed for encoding the message are amortized over all symbols, making it more efficient for long messages. Per symbol, the overhead is only $2/n$ bits, compared to 1 bit per symbol for Huffman coding.

Notably, arithmetic coding implements a first-in-first-out data structure, that is, a queue. Symbols which are encoded first are also decoded first. As we will see, this makes it a convenient choice for compression with autoregressive models (Section 2.2).

2.1.5 Asymmetric numeral systems

More recent examples of streaming codes are asymmetric numeral systems (ANS) [54]. While arithmetic coding implements a queue data structure, ANS operates like a stack. In arithmetic coding, the symbol encoded first is also decoded first (first-in-first-out) while ANS decodes the last symbol first (last-in-first-out). A detailed description is beyond our scope, and we follow Bamler [15] in our summary of the main idea (see Townsend [191] for a different presentation).

Numeral systems like the decimal or binary system can be interpreted as optimal codes for uniform distributions over a finite alphabet (the “base” of the numeral system). They encode a sequence of enumerated symbols into a single integer number, the “stack”. To encode a symbol, we multiply the stack with the base (e.g., 2 or 10) and add the symbol. To decode, we recover the symbol as the stack modulo its base, while we shorten the stack by dividing it by the base. The stack’s length in binary representation is approximately the number of symbols times the logarithm of the base of the numeral system, consistent with entropy coding. Interestingly, the stack’s base can be changed from symbol to symbol and it is still possible to recover the sequence as long as the inverted order of bases is used upon decoding.

ANS generalizes numeral systems from uniform probability distributions to non-uniform ones. Similar to arithmetic coding, symbols are first represented as subintervals of the unit interval. In addition, one discretizes the unit interval using a fine discretization grid. Every point on the grid belongs to the subinterval of a symbol and can be used to represent that symbol. The discretization points form a new alphabet of symbols that can be encoded using a numeral system, as described above. Since there are multiple discretization points in each subinterval, naively encoding one of them leads to a redundant code. However, ANS is able to avoid this redundancy. Bamler [15] notes that the mechanism by which ANS achieves this can be interpreted as a bits-back coding procedure, which will be discussed in Section 2.3.2. Like arithmetic coding, ANS incurs an overhead of up to 2 bits per message.

2.1.6 Continuous models for discrete data

In this subsection, we show how a model for discrete data can be defined and parameterized by a continuous model, a common technique in generative modeling and neural compression. Lossless compression generally operates on discrete data. After all, it is infeasible to losslessly encode any real number with a finite number of bits. However, many generative models in machine learning assume continuous data, and model it with a density function. It turns out such continuous models are still useful for lossless compression.

Suppose we have a density model q over \mathbb{R}^D , and $\mathbf{x} \in \{0, \dots, 255\}^D$ is an RGB image following the ground truth image distribution P . We can derive a PMF over integers by integrating the density over hypercubes, as follows,

$$Q(\mathbf{x}) := \int_{[-.5, .5]^D} q(\mathbf{x} + \mathbf{u}) d\mathbf{u}. \quad (2.5)$$

Assume we add uniform noise $\mathbf{u} \in [-0.5, 0.5]^D$ to the data and that the resulting noisy/dequantized data, $\mathbf{y} = \mathbf{x} + \mathbf{u}$, has density p . Then it is not difficult to show that [183]

$$-\int p(\mathbf{y}) \log_2 q(\mathbf{y}) d\mathbf{y} \geq -\sum_{\mathbf{x}} P(\mathbf{x}) \log_2 Q(\mathbf{x}). \quad (2.6)$$

The left-hand side is the negative log-likelihood which we would optimize when fitting our generative model to the dequantized data. Thus, we can minimize an upper bound on the lossless compression cost under the discretized model (RHS) by fitting a density via maximum likelihood (LHS), provided the discrete data is *dequantized* appropriately with noise [197].

The form of Q as defined in Eq. 2.5, which we call a *discretized density model*, is also useful in itself as a flexible model for discrete data. Examples include PixelCNN++ [160], the prior distribution in a discrete flow [87], as well as entropy models for neural compression [125][14][131].

The integral in Eq. 2.5 in general is intractable to compute, but it can often be broken down into a series of univariate integrals. Therefore let us consider a univariate version of the discretized density model,

$$Q_\theta(x) := \int_{[-0.5, 0.5]} q_\theta(x+u) du.$$

Let F_θ denote the CDF of q_θ , then we can equivalently express the above in terms of a difference of CDF evaluations:

$$Q_\theta(x) := F_\theta(x+0.5) - F_\theta(x-0.5). \quad (2.7)$$

For example, if q_θ is the density of a logistic distribution, then F_θ is the logistic sigmoid function common in deep learning.

2.2 Autoregressive models

Autoregressive models exploit the fact that we can write any probability distribution as a product of conditional distributions using the chain rule of probabilities [23],

$$p(\mathbf{x}) = \prod_i p(x_i | \mathbf{x}_{<i}). \quad (2.8)$$

Here, x_i is the i -th entry of the vector \mathbf{x} and $\mathbf{x}_{<i}$ corresponds to all previous entries in an arbitrary order. The autoregressive factorization does not make any assumptions about the data distribution yet still allows us to easily incorporate useful assumptions. For example, a Markov assumption can be implemented by only considering the entries in $\mathbf{x}_{<i}$ which are close to i . A stationarity assumption is easily incorporated by using the same conditional distribution $p(x_i | \mathbf{x}_{<i})$ at every location. These two assumptions are often reasonable and can drastically reduce the amount of parameters of a model.

Autoregressive modeling lends itself to lossless compression in combination with arithmetic coding (Section 2.1.4) since both deal with data sequentially, one symbol at a time. Each symbol is encoded using the conditional distribution given the data that has already been encoded. This is in contrast to Markov random fields, for example, where conditional distributions are typically only tractable when conditioning on a larger neighborhood. Entropy coding generally requires the number of symbols in each encoding step to be manageable. Autoregressive models provide an important step towards practical entropy coding by decomposing a high-dimensional distribution into low-dimensional conditional distributions. Arithmetic coding is a better match for autoregressive

models than ANS since it operates like a queue – symbols encoded earlier will also be decoded earlier, and therefore will be available as input to conditional distributions.

For data modalities such as audio or text, autoregressive models are an obvious choice due to the signal’s sequential nature. Indeed, two of the top performing algorithms in the Large Text Compression Benchmark [120] use recurrent neural networks to predict the next token of a sentence. While `cmix` [102] uses a mixture of long short-term memory networks (LSTMs) [84] and nonparametric models, `nncp` [16] relies solely on neural network models. However, autoregressive models have long also found application in image compression. For example, JPEG [92] encodes the difference between neighboring DC coefficients, $z_{ij}^{\text{DC}} - z_{i(j-1)}^{\text{DC}}$, in a raster-scan order, which can be thought of as implementing a first-order Markov model.

Mixtures of experts [94][88][182] are a class of autoregressive models proven useful for compressing images,

$$p(x_{ij} \mid \mathbf{x}_{<ij}) = \sum_k \underbrace{p(k \mid \mathbf{x}_{<ij})}_{\text{Gates}} \underbrace{p(x_{ij} \mid \mathbf{x}_{<ij}, k)}_{\text{Experts}}. \quad (2.9)$$

The basic idea behind neural network extensions of this approach is to nonlinearly transform the inputs $\mathbf{x}_{<ij}$ before feeding them into the gates and the experts. For instance, RNADE [197] used a fully connected neural network with a single rectified linear layer to transform the inputs $\mathbf{x}_{<ij}$. Other examples of deep autoregressive models include RIDE [180], PixelRNN [141], or PixelCNN++ [160]. More recent autoregressive modeling papers continue to explore different architectures for transforming $\mathbf{x}_{>ij}$. For example, the Image Transformer [143] uses an architecture based on self-attention [198]. PixelSNAIL [34] uses a combination of convolutions and self-attention layers. Glow [100] uses invertible flows.

With the exception of `nncp` [16] and `cmix` [102], all autoregressive models mentioned so far are *static* models. That is, the models’ parameters are trained once and then remain fixed during the entire encoding and decoding process. In contrast, a *dynamic* model updates its parameters during the encoding process based on already encoded data. The decoder is then able to apply the same model updates based on

the data already received. A simple dynamic autoregressive model for images was proposed by Wu *et al.* [204]. Here, the predictors are linear but the predictor’s parameters are chosen dynamically. This is done by treating a larger neighborhood of preceding pixels as training data for the predictor. Meyer and Tischer [128] improved on this idea by weighting training points based on their distance to the pixel whose value we are predicting.

Autoregressive models come with the restriction that decoding is an inherently sequential procedure. Each symbol can only be decoded after all the symbols have been decoded on which its prediction depends. To improve decoding speed, we can restrict the context $\mathbf{x}_{<i}$ to only a small neighborhood around \mathbf{x}_i , as for example in JPEG. Furthermore, the degree of parallelism can be increased by grouping coordinates of the data into blocks and only modeling the conditional dependencies between blocks, while treating data coordinates within each block as conditionally independent [173][132].

2.3 Latent variable models

Latent variable models represent the data distribution using the sum rule of probability,

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) \, d\mathbf{z} = \int p(\mathbf{x} | \mathbf{z})p(\mathbf{z}) \, d\mathbf{z}, \quad (2.10)$$

where \mathbf{z} is a vector of *latent variables* (also called “latents”), and the joint distribution $p(\mathbf{x}, \mathbf{z})$ factorizes into a *prior* $p(\mathbf{z})$ and a *likelihood* $p(\mathbf{x} | \mathbf{z})$. Latent variable models are ubiquitous in machine learning and include hidden Markov models, mixture models [23], and more recently variational autoencoders (VAEs) [99][150]. By integrating or summing over all possible realizations of the latent vector, latent variable models can capture complex dependencies in the data even when the prior and likelihood take simple forms.

Training a latent variable model by (approximate) maximum-likelihood (see Eq. 2.3) comes with a challenge: unlike in a fully-observed (e.g., autoregressive) model where the data probability $p(\mathbf{x})$ can be readily evaluated, doing so is no longer straightforward since it is

defined through an often intractable integral (Eq. 2.10). Variational inference deals with exactly these complications, and we refer to Zhang *et al.* [217] for more details on these methods. Here, we instead focus on the compression problem, which faces a similar challenge and uses similar tools. Given a latent variable model consisting of a prior $p(\mathbf{z})$ and likelihood $p(\mathbf{x}|\mathbf{z})$, our goal is to compress a given data vector \mathbf{x} with a code length close to its information content under the model, $-\log_2 p(\mathbf{x})$. To simplify the discussion below, we assume discrete \mathbf{z} unless specified otherwise.

2.3.1 Two-part code

We start by considering a simple although generally suboptimal *two-part code* [71]. Here, the data \mathbf{x} is transmitted along with a latent vector \mathbf{z} in two steps. First, the sender decides on a (ideally informative) vector \mathbf{z} , then encodes and transmits it under the prior $p(\mathbf{z})$. Next, \mathbf{x} is compressed and transmitted under the likelihood model, $p(\mathbf{x}|\mathbf{z})$. The receiver then decodes \mathbf{z} , and finally \mathbf{x} given \mathbf{z} , using the same models. Both the prior and likelihood models are often fully factorized, allowing for coding each dimension of \mathbf{z} or \mathbf{x} in parallel, but can also be autoregressive models used in conjunction with arithmetic coding. Assuming negligible entropy coding overhead, the combined code length is

$$l(\mathbf{x}, \mathbf{z}) = -\log_2 p(\mathbf{z}) - \log_2 p(\mathbf{x} | \mathbf{z}) = -\log_2 p(\mathbf{x}, \mathbf{z}). \quad (2.11)$$

Moreover, the sender can (at least in principle) minimize this quantity over all choices of \mathbf{z} , resulting in the code length

$$l_{\text{MTP}}(\mathbf{x}) = \min_{\mathbf{z}} (-\log_2 p(\mathbf{z}) - \log_2 p(\mathbf{x} | \mathbf{z})). \quad (2.12)$$

The minimal two-part code length, $l_{\text{MTP}}(\mathbf{x})$, is generally still suboptimal² [59][200][85]. To see this, consider the following bound on the information content $-\log_2 p(\mathbf{x})$,

²However, it is possible to directly optimize a neural compression method for the two-part code. E.g., Mentzer *et al.* [125] trained a latent variable model with an inference network to minimize the expected coding cost of the two-part code, demonstrating much faster decoding speed than autoregressive models, although at $\sim 20\%$ worse bit-rate.

$$-\log_2 p(\mathbf{x}) \leq -\log_2 p(\mathbf{x}) + d_{\text{KL}}[q(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z} | \mathbf{x})] \quad (2.13)$$

$$= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \mathbf{x})} [-\log_2 p(\mathbf{z}, \mathbf{x})] - H[q(\mathbf{z} | \mathbf{x})], \quad (2.14)$$

also known as the *negative evidence lower bound* (NELBO) of variational inference [217]. Here, q is any distribution over \mathbf{z} which may depend on \mathbf{x} . Crucially, the minimal two-part code length (Eq. 2.12) is a special case of the NELBO where

$$q(\mathbf{z} | \mathbf{x}) = \delta_{\mathbf{z}_{\min}}(\mathbf{z}) \quad (2.15)$$

is a degenerate distribution centered on \mathbf{z}_{\min} , which achieves the minimum in Eq. 2.12. In this case the entropy term vanishes. More generally, the NELBO is minimized when $q(\mathbf{z} | \mathbf{x})$ equals the Bayesian *posterior* distribution $p(\mathbf{z} | \mathbf{x})$, the target of approximate inference [217].

One may naturally wonder whether a coding cost equal to the NELBO can be realized for *any choice of* $q(\mathbf{z} | \mathbf{x})$. Bits-back coding, to be discussed next, answers this question in the affirmative. It further offers a compression interpretation of variational inference: for a given latent variable model, reducing the KL divergence between the distribution q and the posterior distribution is equivalent to minimizing the code length of \mathbf{x} under bits-back coding.

2.3.2 Bits-back coding

While the optimal two-part code chooses a best latent vector \mathbf{z} using deterministic optimization, bits-back coding [60][200][81][85] instead uses a *stochastic* latent code in the form of a sample \mathbf{z} from a distribution q . Surprisingly, this stochastic code can save bits compared to a deterministic code by allowing auxiliary information to piggyback on the choice of the latent code. To gain some intuition, consider again the optimization in Eq. 2.12 for picking the optimal two-part code. When the given data \mathbf{x} is uninformative of what latent variable \mathbf{z} generated it, the minimization over \mathbf{z} may find multiple candidates that are nearly optimal. In the extreme case, ties need to be broken to settle on one of many equally good candidates. Instead of choosing a \mathbf{z}_i randomly, the encoder can do so intentionally so as to communicate additional information via the chosen index i , provided that the decoder can reverse-engineer the encoding procedure.

Before proceeding, we emphasize a connection between sampling and decoding random bits. Given an entropy code derived from a discrete distribution, using it to *decode* a sequence of uniformly random bits produces a random sample from this distribution. For example, consider decoding a random bit-string according to a Huffman tree (Figure 2.1). As we traverse from the root to any leaf node, every edge along the path corresponds to a decision according to a coin flip. The probability of ending up with a particular symbol \mathbf{x} with code length $|C(\mathbf{x})|$ is $2^{-|C(\mathbf{x})|}$, or approximately $2^{\log_2 p(\mathbf{x})} = p(\mathbf{x})$.

In bits-back coding, the sender first generates a sample $\mathbf{z} \sim q(\mathbf{z} | \mathbf{x})$ by *decoding* a random bit-string ξ . Next, the sender encodes \mathbf{x} under $p(\mathbf{x} | \mathbf{z})$, then \mathbf{z} under $p(\mathbf{z})$, and then transmits the resulting bits. The receiver begins by decoding \mathbf{z} and \mathbf{x} from the received bits just as in the two-part code, but then proceeds to recover the exact bit-string the sender used to generate \mathbf{z} . This can be done by *encoding* \mathbf{z} under the distribution $q(\mathbf{z} | \mathbf{x})$, which we assume the receiver has access to. Thus, $-\log_2 q(\mathbf{z} | \mathbf{x})$ extra bits of information are recovered by the decoder (in addition to \mathbf{x}) “free of charge”, giving the scheme its name “bits-back” coding. Finally, and crucially, we note that the initial bits used to generate \mathbf{z} do not have to be truly random. Indeed, to be useful, they should be supplied from a well-compressed bit stream of *auxiliary information* that also needs to be transmitted. For example, if our application requires us to transmit an image along with a thumbnail version of it (e.g., for fast preview), the initial bits can come from the bit-string of the thumbnail image encoded with JPEG. We will examine this issue in more detail shortly.

Figure 2.2 graphically summarizes the algorithm. Given a bit-string ξ of auxiliary bits and the data \mathbf{x} to transmit, the sender starts by decoding a portion of ξ into a stochastic latent code \mathbf{z} of the data, ultimately sending a total number of bits equal to $|\xi| - \log_2 p(\mathbf{z}) - \log_2 p(\mathbf{x} | \mathbf{z}) + \log_2 q(\mathbf{z} | \mathbf{x})$ to communicate both ξ and \mathbf{x} . The net number of bits used for transmitting \mathbf{x} alone is therefore

$$-\log_2 p(\mathbf{z}) - \log_2 p(\mathbf{x} | \mathbf{z}) - (-\log_2 q(\mathbf{z} | \mathbf{x})) \quad (2.16)$$

which corresponds to a discount of $-\log_2 q(\mathbf{z} | \mathbf{x})$ bits compared to a two-part code. Since \mathbf{z} was drawn from $q(\mathbf{z} | \mathbf{x})$, the expected coding cost is exactly the NELBO [217].

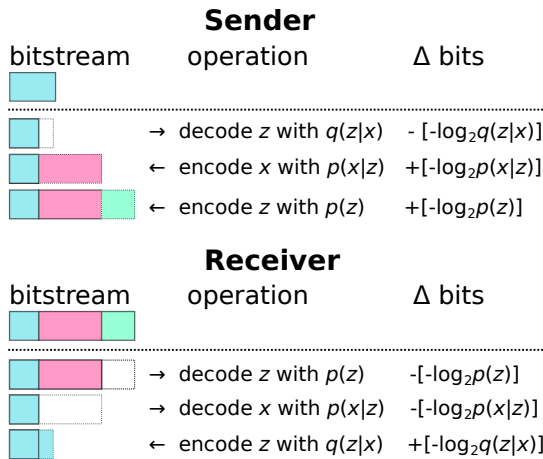


Figure 2.2: An illustration of the encoding (sender) and decoding (receiver) operations of bits-back coding. Auxiliary bits (ξ) are colored in blue.

The choice of auxiliary information is an important one in bits-back coding. If the sender only wishes to communicate a single data sample \mathbf{x} and no additional auxiliary information, the sender may set ξ to be an artificial sequence of random bits, which can still be recovered by the receiver, but contains no useful information. In this case we do not get “bits back”, and suffer an expected code length that is worse than the best two-part code. Even when there are auxiliary bits to transmit, the sender needs to commit to sending at least $-\log q(\mathbf{z}|\mathbf{x})$ many of them for bits-back coding to achieve its nominal efficiency. This is known as the *initial bits problem* [59][193]. One potential solution is to compress a small portion of the data (e.g., a part of an image) with another codec, and use the resulting compressed bit-string for auxiliary bits [192]. If multiple data samples need to be compressed, an elegant solution is *chaining* [193] or “bits-back with feedback” [59]. The idea is to use the running bitstream of previously encoded data samples $\mathbf{x}_1, \dots, \mathbf{x}_{i-1}$ as the source of auxiliary bits for encoding the sample \mathbf{x}_i . The initial bits problem becomes less severe with more data samples, and the code length per sample asymptotically converges to the NELBO. Chaining poses a requirement on the entropy coder that data is decoded in the exact opposite order in which it is encoded. Bits-back is therefore

naturally combined with ANS coding which operates in stack order (Section 2.1.1), popularized by the BB-ANS algorithm [193].

2.3.3 Improvements and extensions

Continuous latents. When the latent vector is continuous, bits-back coding can still be applied after discretizing the latent space and associated distributions [118][193]. More finely discretized latents will require more bits to encode, exacerbating the initial bits problem, but allow us to recover an equal amount of bits back. Bits-back coding thus allows us to code continuous latents with arbitrarily high precision.

Extension to models with multiple latents. Although our discussion focused on a model with a single latent tensor, Ho *et al.* [83] and Townsend *et al.* [192] extended bits-back coding to hierarchical latent variable models, leveraging their superior expressiveness (in terms of better NELBO) to improve compression performance. Bit-Swap [83], in particular, places restrictions on the latent hierarchy (e.g., Markovian) to allow recursive bits-back coding, alleviating the initial-bits problem. Ruan *et al.* [157] and Townsend *et al.* [194] developed bits-back schemes for sequential (e.g., time-series) latent variable models.

Iterative inference. As discussed in Section 2.3.1, the overhead of bits-back coding is equal to the KL divergence between $q(\mathbf{z} | \mathbf{x})$ and the true posterior $p(\mathbf{z} | \mathbf{x})$. The parameters of q (e.g., mean of a Gaussian) are typically predicted from \mathbf{x} by an *amortized inference* network, as in a VAE (see, e.g., [99]). Although cheap to compute, the resulting q distribution is generally suboptimal [40]. Yang *et al.* [209] took a pre-trained model [131] and proposed to directly minimize the KL gap with respect to the q parameters for each given data point \mathbf{x} , resulting in improved image compression performance.

Extended latent spaces. Analogous to the importance-weighted ELBO [32], importance-weighted NELBO provides a tighter variational upper bound on the ideal code length than the NELBO (Eq. 2.14). Ruan *et al.* [157] and Theis *et al.* [181] proposed bits-back coding schemes that operate in an extended latent space and operationalize the coding cost of the importance-weighted NELBO, and demonstrated improved compression performance.

2.4 Invertible flows and other models

Normalizing Flows. Although *continuous* normalizing flows do not lend themselves directly to lossless compression, local bits-back coding [83] allows lossless compression of finely quantized continuous data, with a code length close to the negative log data density up to a constant dependent on the quantization precision. The same method also allows losslessly compressing discrete data using a flow trained on the dequantization objective (e.g., the LHS of Eq. 2.6). By contrast, *discrete normalizing flows* [87][195] directly model integer-valued data and support lossless compression. Such a flow learns a bijection f to map data to an integer latent space, where a factorized prior is assumed; the given data \mathbf{x} can then be compressed by simply entropy-coding its latent representation $f(\mathbf{x})$, and decompressed using f^{-1} . This can be viewed as bits-back coding with a deterministic q .

Other Models. New types of generative models continue to be explored for lossless compression, aiming to achieve greater modeling power (hence better compression rate), while reducing computational demands. For example, discrete diffusion models [168][86] offer parallel encoding/decoding and improved single-image compression rate. Furthermore, Probabilistic Circuits [144] offer lightweight models and enable exact likelihood evaluation and efficient marginalization. Liu *et al.* [110] recently showed that this property makes them a compelling model for lossless compression. By proposing efficient en- and decoding algorithms, Liu *et al.* achieved quasi-linear time complexity in the data dimension, resulting in an order of magnitude faster (de-)compression speed than available integer discrete flows or bits-back coding implementations.

3

Lossy Compression

In the following we review neural *lossy* compression, that is, compression with imperfect data reconstruction. When working with continuous data, such as an analogue signal, lossy compression is necessary, as it is impossible to losslessly store real values with a finite number of bits. In other cases, a lossy reconstruction is often also sufficient, and the flexibility to discard “irrelevant” information allows us to achieve a far lower bit-rate than possible with lossless compression. Lossy compression algorithms for digital media, such as audio, images, and video, routinely obtain an order of magnitude or more bit-rate savings than their lossless counterparts, without noticeable loss of quality to the end users.

We begin in Section 3.1 with basic background on lossy compression in the rate-distortion setting, reviewing the classical rate-distortion theory and algorithms such as vector quantization and transform coding. The rate-distortion VAE is introduced as a conceptual lossy compressor, embodying many approaches to be discussed. Section 3.2 introduces neural lossy compression as a natural extension of transform coding, and discusses high-level architecture and modeling choices, with connections to hierarchical latent variable models. We then discuss the core issue of end-to-end learning quantized representations and entropy modeling

in Section 3.3. In Section 3.4, we introduce an emerging technique of compression without quantization, which can be seen as operationalizing the bit-rate cost of a rate-distortion VAE, as given by a KL divergence. In aiming towards perceptually-pleasant reconstructions of data such as images, in Section 3.5 we discuss various techniques for evaluating perceptual quality of images, and show how neural lossy compression can be optimized for high perceptual quality, in addition to the rate-distortion trade-off. We further consider the setting where the distortion is defined according to performance on downstream tasks in Section 3.6, and finally give a brief overview of neural video compression in Section 3.7.

Frequently, we make reference to a *quantization* operation, denoted by $\llbracket \cdot \rrbracket$. Quantization can be implemented in various ways in neural compression (e.g., by rounding to the nearest integer, denoted by $\lfloor \cdot \rfloor$; we give a detailed treatment in Section 3.3), but at a high level is understood to be any mapping, either deterministic or stochastic, from a (often continuous) set of input values to a smaller and *countable* set of output values.

3.1 Background

We give a brief overview of the main results of lossy compression, in the classical rate-distortion setting studied by Shannon [164]. Here, lossy compression involves encoding the given data into a *discrete* representation, which is communicated losslessly to the receiver, using as few bits as possible. The quality of lossy reconstructions is measured by a fidelity criterion, or a *distortion* function between pairs of original data samples and reconstructions. We refer readers to surveys by Berger [19] and Gray and Heuhoff [69] for more comprehensive coverage of classical results and historical developments in lossy compression.

3.1.1 Rate-distortion theory

For our purpose, we formalize a lossy compression algorithm (or a lossy *codec*) as a 3-tuple consisting of an encoder, a decoder, and an entropy code, denoted by $c = (\mathbf{e}, \mathbf{d}, \gamma)$. The *encoder* $\mathbf{e} : \mathcal{X} \rightarrow \mathcal{W}$ maps each data point $\mathbf{x} \in \mathcal{X}$ to a representation \mathbf{w} in a countable set \mathcal{W} (this

can be the set of natural numbers, without loss of generality). The decoder $\mathbf{d} : \mathcal{W} \rightarrow \hat{\mathcal{X}}$ maps each representation \mathbf{w} to a *reconstruction point* $\hat{\mathbf{x}} \in \hat{\mathcal{X}}$.

The encoder and decoder further agree to transmit the symbols of \mathcal{W} losslessly with an *entropy code* γ (see Section 2.1.1); we write $l(\mathbf{x}) := |\gamma(\mathbf{e}(\mathbf{x}))|$ to denote the code length assigned to the encoding of \mathbf{x} . Suppose we are given a *distortion function* $\rho : \mathcal{X} \times \hat{\mathcal{X}} \rightarrow [0, \infty)$ that measures the error caused by representing \mathbf{x} with the lossy reconstruction $\hat{\mathbf{x}}$. Historically, this is typically some form of a squared error, i.e., $\rho(\mathbf{x}, \hat{\mathbf{x}}) \propto \|\mathbf{x} - \hat{\mathbf{x}}\|^2$, and we will revisit this choice and discuss alternatives in Section 3.5. Given the above, lossy compression is then generally concerned with minimizing the average *distortion*

$$\mathcal{D} = \mathbb{E}_{\mathbf{x} \sim P_{data}} [\rho(\mathbf{x}, \hat{\mathbf{x}})], \quad (3.1)$$

where $\hat{\mathbf{x}} := \mathbf{d}(\mathbf{e}(\mathbf{x}))$, and simultaneously, the (*bit-*)*rate*

$$\mathcal{R} = \mathbb{E}_{\mathbf{x} \sim P_{data}} [l(\mathbf{x})], \quad (3.2)$$

that is, the average number of bits needed to encode the data. We want to minimize these two quantities with respect to the choice of the encoding and decoding procedures \mathbf{e} and \mathbf{d} , as well as the entropy code γ .

Rate-distortion (R-D) theory [164][39] establishes limits on the performance of any such algorithm. Any lossy compression algorithm implements a noisy channel that receives a data input \mathbf{X} and outputs a reconstruction $\hat{\mathbf{X}}$, described by a conditional distribution $Q_{\hat{\mathbf{X}}|\mathbf{X}}$. The mutual information, $I[\mathbf{X}, \hat{\mathbf{X}}]$, is then defined as the KL divergence between the joint distribution $P_{data} \cdot Q_{\hat{\mathbf{X}}|\mathbf{X}}$ and the product of its two marginal distributions. As we will also see in Section 3.4, the mutual information $I[\mathbf{X}, \hat{\mathbf{X}}]$ is a fundamental measure of the amount of information needed to transmit $\hat{\mathbf{X}}$ given \mathbf{X} (in bits per sample), and for a given distortion threshold D , the lowest achievable bit-rate is characterized by the information *rate-distortion function* [164],

$$\mathcal{R}_I(D) = \inf_{Q_{\hat{\mathbf{X}}|\mathbf{X}}: \mathbb{E}[\rho(\mathbf{X}, \hat{\mathbf{X}})] \leq D} I[\mathbf{X}, \hat{\mathbf{X}}]. \quad (3.3)$$

The (information) R-D function is a fundamental quantity that depends only on the data distribution and choice of distortion function, and

generalizes the Shannon entropy H (Eq. 2.1) from lossless compression to lossy compression. The R-D function has a few general properties, e.g., it is always non-increasing and convex, but is otherwise unknown analytically outside of a few cases. It can be numerically estimated by the Blahut-Arimoto algorithm [24][9] when the data and reconstructions are discrete and low-dimensional, and recent work has tackled the more general setting and estimated the R-D functions of continuous and high-dimensional data such as natural images [211].

In theory, the optimal rate $\mathcal{R}_I(D)$ is achievable by vector quantization [39], to be described in Section 3.1.3, but this approach to compression quickly becomes intractable for high-dimensional data. Rather than trying to achieve the theoretically optimal rate $\mathcal{R}_I(D)$ with any codec at any computational cost, in practice the design of a lossy codec is constrained by practical considerations, such as the computation budget of the target hardware, or the decoding latency acceptable for the application. Denoting the set of all the acceptable codecs under consideration by \mathcal{C} , we can instead consider an *operational* rate-distortion function¹, formalized by

$$\mathcal{R}_O(D) = \inf_{c \in \mathcal{C}: \mathbb{E}[\rho(\mathbf{X}, \hat{\mathbf{X}})] \leq D} \mathbb{E}[l(\mathbf{X})]. \quad (3.4)$$

Compared to Eq. 3.3, we replaced mutual information by the operational rate and optimize over an actual lossy codec c . We can relax the constrained optimization problem to an unconstrained one, by introducing the rate-distortion Lagrangian [24][37],

$$L(\lambda, c) = \mathcal{R}(c) + \lambda \mathcal{D}(c) = \mathbb{E}[l(\mathbf{X})] + \lambda \mathbb{E}[\rho(\mathbf{X}, \hat{\mathbf{X}})]. \quad (3.5)$$

For each fixed $\lambda > 0$, the minimum of this objective yields a codec c^* whose operational distortion-rate performance, $(\mathcal{D}(c^*), \mathcal{R}(c^*))$, lies on

¹In usual treatments of R-D theory (see e.g., [39]), *the* operational R-D function is defined as the asymptotic rate achievable by compressing multiple samples jointly using *any* lossy codec, in the limit of infinitely many jointly compressed samples. This is different from our definition here, which instead focuses on the operational R-D performance achievable within a constrained family of codecs \mathcal{C} , by compressing a single sample at a time (as is common in applications). Thus we generally only have $\mathcal{R}_O(D) \geq \mathcal{R}_I(D)$, i.e., $\mathcal{R}_I(D)$ is in general no longer achievable in this setup.

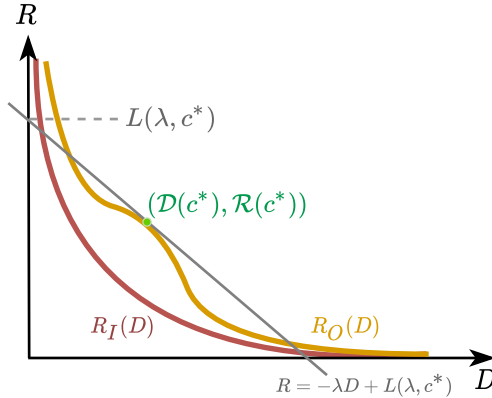


Figure 3.1: Visualizing the operational R-D optimization problem. We adjust our codec c to minimize the R -axis intercept of a straight line (gray) with slope $-\lambda$ and passing through $(\mathcal{D}(c), \mathcal{R}(c))$; an optimum occurs when the line becomes tangent to the operational R-D curve \mathcal{R}_O at the point $(\mathcal{D}(c^*), \mathcal{R}(c^*))$. Note that the operational R-D curve (orange) is not necessarily convex, and always lies above the information R-D curve \mathcal{R}_I (red), i.e., $\mathcal{R}_O(D) \geq \mathcal{R}_I(D), \forall D$.

the convex hull of the operational R-D curve². This is illustrated in Figure 3.1. The hyperparameter λ can be interpreted as a Lagrange multiplier, and codecs with different operational rate-distortion trade-offs can be found by minimizing the Lagrangian with various λ . Most current end-to-end learned lossy compression methods to be discussed in Section 3.2 follow this approach, training one codec for each λ . We note that it is not always possible to attain every point on the operational R-D curve with this approach [43], and alternative methods based on constrained optimization have been proposed [145][156].

3.1.2 Connection to latent variable modeling and R-D VAEs

In Section 2, we saw that an integral part of lossless compression is estimating a good model of the data. In lossy compression, a similar connection can be drawn between rate-distortion optimization and data modeling, via a particular class of latent variable models [12][184][6][211] which we will term *rate-distortion variational autoencoders*.

²In practice, due to non-convexity, we typically reach a local minimum of the Lagrangian, and therefore an R-D point lying above the operational R-D curve.

Let \mathcal{Z} be any latent space, and $g : \mathcal{Z} \rightarrow \hat{\mathcal{X}}$ any measurable function. Yang and Mandt [211] showed that an upper bound on the rate-distortion function $\mathcal{R}_I(D)$ can be obtained by optimizing the following objective:

$$\begin{aligned} \mathcal{L}(q(\mathbf{z}|\mathbf{x}), p(\mathbf{z}), g, \lambda) &:= \mathbb{E}_{\mathbf{x} \sim P_{data}} [d_{\text{KL}}[q(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z})]] \\ &\quad + \lambda \mathbb{E}_{\mathbf{x} \sim P_{data}, \mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\rho(\mathbf{x}, g(\mathbf{z}))], \end{aligned} \quad (3.6)$$

where $q(\mathbf{z}|\mathbf{x})$ and $p(\mathbf{z})$ are arbitrary distributions defined on \mathcal{Z} . The objective comes from a Lagrangian relaxation of the constrained optimization problem defining $\mathcal{R}_I(D)$ (in Eq. 3.3), with the expected KL divergence term serving as a variational upper bound on $I[\mathbf{X}, \hat{\mathbf{X}}]$. Suppose there exists a conditional density *aligned with* the given distortion function ρ , in the sense that

$$p(\mathbf{x}|\mathbf{z}) = C \exp\{-\rho(\mathbf{x}, g(\mathbf{z}))\},$$

where the normalizing constant C does not depend on \mathbf{z} . In the common case of a squared error distortion, $p(\mathbf{x}|\mathbf{z})$ is an isotropic Gaussian density with mean $g(\mathbf{z})$ and constant variance. Then, the Lagrangian \mathcal{L} can be easily seen to equal the NELBO objective of a (β -) VAE (up to a constant) [12][184][211],

$$\begin{aligned} \mathcal{L}(q(\mathbf{z}|\mathbf{x}), p(\mathbf{z}), g, \lambda) &= \mathbb{E}_{\mathbf{x} \sim P_{data}} [d_{\text{KL}}[q(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z})]] \\ &\quad + \lambda \mathbb{E}_{\mathbf{x} \sim P_{data}, \mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [-\log p(\mathbf{x}|\mathbf{z})] + \text{const}, \end{aligned} \quad (3.7)$$

in which $p(\mathbf{x}|\mathbf{z})$ is the likelihood/observation model, $q(\mathbf{z}|\mathbf{x})$ is the variational posterior, $p(\mathbf{z})$ is the (variational) prior, and g is the decoder network. To make the correspondence with a VAE exact (“ $\beta = 1$ ”), we can also absorb λ into the definition of the likelihood model, positing a density $p(\mathbf{x}|\mathbf{z}) = C \exp\{-\lambda \rho(\mathbf{x}, g(\mathbf{z}))\}$. In the case of a squared distortion, this corresponds to a Gaussian density with a precision (inverse variance) parameter proportional to λ . Given the above equivalence between Eq. 3.6 and Eq. 3.7, we refer to the model associated with Eq. 3.6 as a *rate-distortion VAE* (or *R-D VAE*), even when there may not exist a proper likelihood density $p(\mathbf{x}|\mathbf{z})$ aligned with the given distortion function ρ .

Besides yielding computationally tractable bounds on the R-D function [211], rate-distortion VAEs also provide a useful perspective on

many of the lossy compression algorithms we will encounter. Most existing neural lossy codecs can be seen as instances of R-D VAEs with a discrete latent space (Section 3.2.2), and many switch to a continuous latent space for end-to-end training [12] (Section 3.3), while methods that bypass quantization (Section 3.4) almost always use a continuous latent space [2][186].

3.1.3 Vector Quantization

Vector quantization (VQ), a classical technique from signal processing, is perhaps the most basic and general form of a lossy codec. A vector quantization scheme, or a vector quantizer, consists of a set of integer representations $\mathcal{W} = \{1, 2, \dots, k\}$, an encoder \mathbf{e} assigning a given data to an integer representation $w \in \mathcal{W}$, and a decoder \mathbf{d} that returns a reconstruction given w , usually implemented by indexing into a codebook of k reconstruction vectors. Note that there is no restriction on the encoding and decoding functions, other than the cardinality of \mathcal{W} . The goal is then to determine an optimal quantization scheme for a given data distribution, under some objective. Commonly, the objective is to minimize a reconstruction error (Eq. 3.1), as in the k -means algorithm, but can also more generally be an operational rate-distortion trade-off (Eq. 3.4). An optimal quantizer can be shown to always encode any given data point to its “nearest-neighbor” in the codebook to minimize the reconstruction error ρ (or, a combination of reconstruction error and code length, in the rate-constrained version), and set the j th entry in the codebook $\hat{\mathbf{x}}_j$ to minimize the expected reconstruction error between $\hat{\mathbf{x}}_j$ and data points assigned to index j [37].

For some data source, e.g., the uniform or the Laplace distribution, optimal quantization can be characterized analytically [223][176]. In most applications, Lloyd-Max-style algorithms (more widely known as k -means in machine learning) [114][37] are instead used to estimate a quantization scheme from data samples. This usually involves minimizing an empirical rate-distortion cost (Eq. 3.5) over a dataset, which is still a basic ingredient in today’s learned compression approaches [10] (see Section 3.2).

Given unlimited data and compute, VQ can approximate any data distribution arbitrarily well. Indeed, the theoretical limit of lossy compression, the rate-distortion function $\mathcal{R}_I(D)$, can be shown to be achievable by jointly quantizing multiple data samples together with increasingly long blocks [39]. However, VQ comes with the severe downside of poor scalability and data efficiency [63]. The computational and storage demand of VQ increases quickly with the data dimensionality. In high dimensions, an exceedingly large number of quantization points and high amounts of training data are needed to approximate the data distribution well. As a result, VQ is typically found in low-rate applications, such as low-rate speech coding [162]. We note however, k -means-style vector quantization has been successfully applied in a convolutional latent space for image generation [142], and serves as an image tokenizer in many recent transformer-based models for image [56][213] and text-to-image generation [149][214].

3.1.4 Transform coding

Instead of quantizing the data directly, it is often much easier to do so in a transformed space, where the representation of the data is uncorrelated and scalar quantization can be effectively applied. The core idea behind transform coding [67] is therefore to divide the task of lossy compression into decorrelation and quantization: first, the sender applies an *analysis transform* f to data \mathbf{x} , resulting in (ideally decorrelated) transform coefficients, $\mathbf{z} = f(\mathbf{x})$; and second, coordinate-wise scalar quantization $\llbracket \cdot \rrbracket$ is applied to obtain a discretized representation $\hat{\mathbf{z}} = \llbracket \mathbf{z} \rrbracket$. The symbols representing $\hat{\mathbf{z}}$ can then be converted to a bit-string by entropy coding (discussed in Section 2.1.1) under an entropy model $P(\hat{\mathbf{z}})$. The receiver losslessly recovers $\hat{\mathbf{z}}$, and computes a reconstruction $\hat{\mathbf{x}} = g(\hat{\mathbf{z}})$ using a *synthesis transform* g , which is often the inverse of the analysis transform. In the terminology of Section 3.1.1, the encoder $\mathbf{e} = \llbracket \cdot \rrbracket \circ f$ is the composition of the analysis transform f followed by quantization, and the decoder \mathbf{d} is implemented by the synthesis transform g .

The analysis transforms is typically linear and invertible, and is implemented as multiplication with an orthogonal matrix. A classic example is the Karhunen-Loève Transform (KLT), which is the same

as principle component analysis (PCA) in our setting. The KLT maps data samples into the eigen-basis of the covariance matrix of the data distribution, and produces uncorrelated transform coefficients \mathbf{z} . Another example is the discrete cosine transform (DCT), which we saw in JPEG in Section 1. DCT can be shown to approximate PCA asymptotically [5], and is one of the most widely used transforms in signal processing and compression. The choice of an orthogonal transform simplifies the theoretical analysis and design of transform coding algorithms, and for Gaussian data, the KLT followed by uniform scalar quantization can be shown to be optimal [68].

Unlike vector quantization, which effectively optimizes over all possible choices of encoder $\mathbf{e} : \mathcal{X} \rightarrow \mathcal{W}$ and decoder $\mathbf{d} : \mathcal{W} \rightarrow \hat{\mathcal{X}}$, transform coding implicitly restricts the solution space of allowable codecs, e.g., the set of quantization points must be in the range of g , and therefore generally cannot achieve the unconstrained optimal performance of VQ. The lack of theoretical optimality of transform coding is more than made up for by its vastly superior scalability over VQ, as evidenced by its wide-spread use in the compression of digital media, such as images and videos [67].

3.2 Neural lossy compression

Most current neural lossy compression methods are based on the paradigm of *nonlinear transform coding* [10] and use learned functions to encode data into a discrete representation, typically by quantizing a continuous representation. Compared to linear transforms (discussed in Section 3.1.4), non-linear transforms are more flexible and can better adapt to the data distribution [10], can be optimized for custom losses (such as perceptual losses in Section 3.5), and are demonstrated to be optimal in some cases [199]. In the following, we present high-level architectures and modeling choices, deferring a full discussion on model training and entropy modeling to Section 3.3.

3.2.1 Overview

In neural compression with non-linear transform coding, we replace the various components of transform coding (Section 3.1.4), such as the analysis transform f , synthesis transform g , and entropy model P , with neural networks or other function approximators, and learn them end-to-end on the data of interest.

As in transform coding, we obtain a discrete representation $\hat{\mathbf{z}}$ by quantizing the transform coefficients \mathbf{z} , i.e., $\hat{\mathbf{z}} = \llbracket \mathbf{z} \rrbracket$, $\mathbf{z} = f(\mathbf{x})$. It's common to refer to $\hat{\mathbf{z}}$ as a tensor of “latents”, due to connections to latent variable models (see Sections 3.1.2, 3.2.2, 3.3.4), although there is some ambiguity in this terminology and some papers also apply it to the continuous transform coefficients \mathbf{z} . As before, an entropy model P is used to losslessly transmit $\hat{\mathbf{z}}$ via bit-strings (Section 2.1.1). The entropy models are often based on powerful models from neural lossless compression (Section 2). Most commonly, the objective is to simultaneously minimize the rate,

$$\mathcal{R} := \mathbb{E}[-\log_2 P(\llbracket f(\mathbf{X}) \rrbracket)], \quad (3.8)$$

and the distortion,

$$\mathcal{D} := \mathbb{E}[\rho(\mathbf{X}, g(\llbracket f(\mathbf{X}) \rrbracket))].$$

The expectations are taken w.r.t. the data distribution P_{data} and are estimated with data samples. Unlike in Section 3.1.1, here we no longer concern ourselves with an entropy code γ . Instead, we directly optimize the information content in place of a code length in Eq. 3.8, knowing that an entropy code can (in principle) always be derived from P such that $-\log_2 P(\cdot) \approx |\gamma(\cdot)|$ (see Section 2.1.1). If we denote the true marginal distribution of $\hat{\mathbf{Z}}$ by P^* (which is induced by P_{data} and the encoding procedure, via $\hat{\mathbf{Z}} = \llbracket f(\mathbf{X}) \rrbracket$), then the rate loss can be equivalently written as the cross entropy,

$$\mathcal{R} = \mathbb{E}_{\hat{\mathbf{z}} \sim P^*(\hat{\mathbf{z}})}[-\log_2 P(\hat{\mathbf{z}})], \quad (3.9)$$

which precisely captures the cost of entropy coding $\hat{\mathbf{z}} \sim P^*(\hat{\mathbf{z}})$ under our model $P(\hat{\mathbf{z}})$, and serves as an upper bound to the entropy $H[\hat{\mathbf{Z}}]$.

As in Eq. 3.5, we would like to optimize the operational R-D performance in the form of a rate-distortion Lagrangian,

$$\mathbb{E}[-\log_2 P(\llbracket f(\mathbf{X}) \rrbracket)] + \lambda \mathbb{E}[\rho(\mathbf{X}, g(\llbracket f(\mathbf{X}) \rrbracket))], \quad (3.10)$$

for a suitable choice of the trade-off factor λ . However, as is written, the objective is not suitable for end-to-end training with SGD, due to the non-differentiability of the quantization operator and the rate loss.

We dedicate Section 3.3 to a detailed discussion of the various approaches to quantization, differentiable approximations, and the related topic of entropy modeling. In the rest of this section, we will examine other aspects and design choices of non-linear transform coding,

3.2.2 Connection to variational autoencoders

The R-D objective in Eq. 3.10 closely resembles that of an autoencoder [29][158], where the analysis and synthesis transforms f and g correspond to the encoder³ and decoder of an autoencoder. The rate term can be seen as a regularizer imposed on the representation $\hat{\mathbf{z}}$, and differs from the regularization in traditional autoencoders based on dimensionality reduction or sparsity [66].

Another way to interpret such a (deterministic) autoencoder is by viewing it as a *variational* autoencoder, but with a degenerate variational posterior distribution. Let us consider a rate-distortion VAE (introduced in Section 3.1.2) with discrete latent variables $\hat{\mathbf{z}}$ ⁴, variational posterior $q(\hat{\mathbf{z}}|\mathbf{x})$ and prior $P(\hat{\mathbf{z}})$. For each fixed \mathbf{x} , if we define $q(\hat{\mathbf{z}}|\mathbf{x}) := \delta_{\llbracket f(\mathbf{x}) \rrbracket}(\hat{\mathbf{z}})$, i.e., we let the approximate posterior concentrate all its mass on the quantized transform coefficients $\llbracket f(\mathbf{x}) \rrbracket$, then the R-D cost of non-linear transform coding in Eq. 3.10 is equal to the NEBLO in Eq. 3.7. Specifically, given a data sample \mathbf{x} , the bit-rate of the R-D VAE reduces to the code length in Eq. 3.10:

³Unfortunately, this use of the term “encoder” for the function f of an autoencoder clashes with our definition of encoders \mathbf{e} in Section 3. The “encoder” can also mean more abstractly the party initiating the data communication (and similarly, “decoder” can refer to the party receiving the data). The meaning is usually clear from the context.

⁴The notation for latent variables in Section 3.1.2 was “ \mathbf{z} ” instead. Unfortunately this collides with the use of “ \mathbf{z} ” as (unquantized) transform coefficients in this section.

$$d_{\text{KL}}[q(\hat{\mathbf{z}}|\mathbf{x})\|P(\hat{\mathbf{z}})] = -\log P(\llbracket f(\mathbf{x}) \rrbracket). \quad (3.11)$$

This situation is analogous to the one in Section 2.3.1, where the two-part code can be considered a special case of bits-back coding with a deterministic posterior distribution. We achieve a coding cost equal to the KL-divergence (LHS of the above equation) by simply entropy coding the deterministic configuration of q under model P .

Do other kinds of R-D VAEs, especially ones with non-degenerate, or stochastic posterior distributions, also have a role to play in lossy compression algorithms? We will address this question in Section 3.3.4 and Section 3.4.

3.2.3 Neural Transform Architectures

Feedforward neural networks are most often used for the encoding and decoding transforms (f, g) in lossy compression, as in traditional autoencoders. For compressing unstructured data, fully-connected neural networks have been used [10][211]. In image compression, the networks are typically convolutional neural networks (CNNs), with f implementing downsampled convolutions, and g typically implementing upsampled convolutions [11] or sub-pixel convolutions [184][189]. Various architectural improvements have been made to better capture and remove redundancies in the input data; these include replacing the usual ReLU activation with Generalized Divisive Normalization [11], introducing residual connections [35], attention mechanisms [112][35], and vision transformers [222].

We note that in the low-distortion / high-rate regime, a sufficiently large latent space [211] and transforms with sufficient capacity (e.g., as determined by the number of filters in a CNN) are generally needed to maximize the rate-distortion performance of a given architecture [12][10]. Invertible networks have also been shown to provide good inductive bias for this setting [77][205].

Recurrent and hierarchical architectures. Instead of encoding data \mathbf{x} into \mathbf{z} with a single neural network, it can be beneficial to introduce a feedback mechanism to involve the decoder in the encoding process. We may divide the compression of \mathbf{x} into T stages, each stage generating an incremental representation $\hat{\mathbf{z}}_t$ from the encoder to the

decoder, and the final $\hat{\mathbf{z}}$ consisting of the concatenation $\hat{\mathbf{z}}_1, \hat{\mathbf{z}}_2, \dots, \hat{\mathbf{z}}_T$. At the beginning of stage t , the decoder has available a tentative data reconstruction $\hat{\mathbf{x}}_{t-1}$ computed from stage $t-1$, using the already received $\hat{\mathbf{z}}_1, \hat{\mathbf{z}}_2, \dots, \hat{\mathbf{z}}_{t-1}$. Crucially, the encoder is equipped with a copy of the decoder, so having computed the same $\hat{\mathbf{x}}_{t-1}$, the encoder only encodes the information in \mathbf{x} that is not present in $\hat{\mathbf{x}}_{t-1}$ (e.g., by encoding the *residual*, $\mathbf{x} - \hat{\mathbf{x}}_{t-1}$, in image compression). The resulting representation $\hat{\mathbf{z}}_t$ is sent to the decoder, which then uses $\hat{\mathbf{z}}_t$ to compute an improved tentative reconstruction $\hat{\mathbf{x}}_t$. This process continues, with $\hat{\mathbf{x}}_T$ declared the final data reconstruction $\hat{\mathbf{x}}$. Here, $f : \mathbf{x} \rightarrow \hat{\mathbf{z}}$ (similarly, g) is no longer implemented by a feedforward architecture, but rather a recurrent one, with $\hat{\mathbf{x}}_t$ being the recurrent state.

Such an approach embodies the idea of “analysis-by-synthesis”, an influential model of perception and comprehension [20], and enables progressive compression whereby the data reconstruction $\hat{\mathbf{x}}_t$ improves as more information in $\hat{\mathbf{z}}_t$ is transmitted. Many video compression methods in Section 3.7 also follow the same predictive coding paradigm, with $\mathbf{x} = [\mathbf{x}_1, \dots, \mathbf{x}_T]$ being a sequence of video frames, and $\hat{\mathbf{z}}_t$ carrying the information in the reconstructed frame $\hat{\mathbf{x}}_t$.

We give a concrete example by Toderici *et al.* [188][189], who proposed some of the first recurrent neural architectures for progressive and variable-rate image compression. Here, the computation at each stage can be summarized as

$$\begin{aligned}\hat{\mathbf{z}}_t &= \llbracket f_t(\mathbf{r}_{t-1}) \rrbracket, \mathbf{r}_t = \mathbf{x} - \hat{\mathbf{x}}_t, \hat{\mathbf{x}}_t = g_t(\hat{\mathbf{z}}_t) + \alpha \hat{\mathbf{x}}_{t-1}; \\ \mathbf{r}_0 &= \mathbf{x}, \hat{\mathbf{x}}_0 = \mathbf{0},\end{aligned}$$

where \mathbf{r}_t , f_t , and g_t denote the residual vector, encoding transform, and decoding transform at time t , respectively, and $\alpha \in \{0, 1\}$ allows two different modes of operation. With $\alpha = 1$ (“additive reconstruction” mode), f_t and g_t can simply be a pair of CNNs (with a separate pair for each t), and are trained to additively correct the cumulative reconstruction at each stage. With $\alpha = 0$ (“one-shot reconstruction” mode), f_t and g_t are chosen to be stateful, and typically LSTM architectures, which are trained to directly predict the entire original image at each stage. Progressive and variable-rate image compression can then be achieved by controlling the number of steps T of the recursive computation.

The idea of reconstructing data progressively from abstract to increasingly specific concepts has also been well-established in deep generative modeling, often in the form of hierarchical latent variable models [70] [169] [101]. Many of these generative models provide inspirations for compression. The recurrent architecture for compression discussed above is intimately related to bidirectional inference in hierarchical VAEs [169]. Hierarchical and latent variable modeling is also prevalent in the entropy models of neural lossy codecs, in the form of a hyperprior [14] and its extensions [131], [132]. Yang and Mandt [211] trained a deep ResNet-VAE [101] to estimate bounds on the R-D function (see Section 3.1.2) and the theoretical performance of compression without quantization (to be discussed in Section 3.4), while Duan *et al.* [52] trained a similar model with additive uniform noise (to be discussed in Section 3.3.3) to obtain improvements in practical lossy compression performance.

3.3 Learned quantization and rate control

Neural networks have been used in image compression since before 1990 [170][96], but techniques were only recently developed to allow end-to-end training directly on a rate-distortion objective [188][11][184]. The main stumbling block has been the fact that the quantization operation and the discrete rate loss (see Eq. 3.10) are not differentiable. As quantization maps (usually continuous) input to a discrete set, its derivative is zero almost everywhere and undefined at points of discontinuity. By the chain rule, the parameters of the encoder transform also receive zero gradient almost everywhere. Moreover, the rate loss, defined via the PMF $P(\hat{\mathbf{z}})$, also has no derivative w.r.t. the discrete representation $\hat{\mathbf{z}}$.

Below we survey the major approaches developed over the years for dealing with the non-differentiability problem, organized by how quantization is done. Since the choice of quantization affects the choice of the entropy model, as well as strategies for optimizing the bit-rate of quantized representations $\hat{\mathbf{z}}$ (Eq. 3.9), we discuss these issues jointly.

We note that although most of these techniques are developed for lossy compression, they are equally useful in lossless compression methods that make use of quantization [87], [125].

3.3.1 Binarization

Some of the earliest work in neural lossy compression obtained discrete representations $\hat{\mathbf{z}}$ by binarizing the output of the encoder (analysis) network. The bit-rate was controlled in various ways when training the neural transforms; afterwards, a separate entropy model was fit to the empirical distribution of $\hat{\mathbf{z}}$ and ultimately used for compression. We discuss specific approaches below.

Based on techniques for training binarized neural networks [202], Toderici *et al.* [188][189] proposed stochastic binarization: each scalar element z of \mathbf{z} is separately mapped to the interval $[-1, 1]$ (e.g., using point-wise tanh as the last layer of f), and then stochastically rounded to -1 or 1 based on how close z is to either value:

$$\llbracket z \rrbracket := \text{B}(z) = z + \epsilon, \quad \mathbb{P}(\epsilon) = \begin{cases} \frac{1+z}{2}, & \epsilon = 1 - z \\ \frac{1-z}{2}, & \epsilon = -1 - z \end{cases}$$

To backpropagate through stochastic binarization, they used a Straight-through Estimator (STE) [17] and defined the gradient to be that of the identity function,

$$\frac{d}{dz} \text{B}(z) := \frac{d}{dz} \mathbb{E}[\text{B}(z)] = \frac{d}{dz} z = 1.$$

Toderici *et al.* [188][189] trained the neural transforms (f, g) to only minimize the distortion loss, relying on constraints on the dimensionality of the binary $\hat{\mathbf{z}}$ to implicitly control the rate. After training, a separate autoregressive entropy model (similar to a PixelRNN) is learned on the empirical distribution of $\hat{\mathbf{z}}$ to further reduce the bit-rate [189].

Li *et al.* [108] deterministically binarized \mathbf{z} to $\{0, 1\}$ and used STE for backpropagation, optimizing a surrogate rate-distortion objective. For rate control, they introduced a learned masking mechanism to the encoder network to encourage sparsity in $\hat{\mathbf{z}}$, and optimized a surrogate rate loss defined in terms of the soft count of non-zeroed-out elements of $\hat{\mathbf{z}}$. After training, they fit a separate PixelCNN-style entropy model to further reduce rate, similar to Toderici *et al.* [189].

3.3.2 Soft-to-Hard Vector Quantization

Agustsson *et al.* [1] proposed to use vector quantization with learned codebook values and introduced corresponding techniques for differentiable quantization and rate control. They considered the k -means-style (hard) quantization operation (see Section 3.1.3), mapping \mathbf{z} to its closest codebook vector,

$$\llbracket \mathbf{z} \rrbracket := \text{VQ}(\mathbf{z}, \mathcal{C}) = \mathbf{c}_j, \quad \text{with } j = \arg \min_i \|\mathbf{z} - \mathbf{c}_i\|,$$

where $\mathcal{C} = \{\mathbf{c}_i | i = 1, 2, \dots, k\}$ is a finite set of codebook vectors in \mathbb{R}^n learned alongside the model. Agustsson *et al.* [1] proposed to approximate hard quantization by a differentiable *soft* quantization, via a linear combination of the codebook vectors weighted by how close they are to \mathbf{z} :

$$\text{VQ}(\mathbf{z}, \mathcal{C}) \approx \text{SoftQ}(\mathbf{z}, \mathcal{C}) := \sum_i \phi_i \mathbf{c}_i.$$

Here $\phi \in \Delta^{N-1}$ is a probability vector computed as the softmax of weighted distances, $\phi = \phi(\mathbf{z}, \mathcal{C}) := \text{softmax}(-\sigma[\|\mathbf{z} - \mathbf{c}_1\|^2, \dots, \|\mathbf{z} - \mathbf{c}_M\|^2])$, with $\sigma > 0$ a hyperparameter. In practice, due to the prohibitive computation cost of VQ, the proposed method is only applied independently to small blocks of \mathbf{z} . By annealing $\sigma \rightarrow \infty$ throughout training, soft quantization gradually approaches hard quantization, and a proper annealing schedule is needed to ensure effective training. For rate control, Agustsson *et al.* [1] optimized a surrogate rate loss based on the empirical distribution of the soft assignment probabilities, estimated via histograms.

Mentzer *et al.* [124] simplified the above technique by performing only scalar quantization and dispensing with the annealing procedure. They fixed σ at a constant (usually 1), and applied STE to differentiate through (hard) quantization using the gradient of soft quantization,

$$\frac{\partial}{\partial z} \text{VQ}(z, \mathcal{C}) := \frac{\partial}{\partial z} \text{SoftQ}(z, \mathcal{C}).$$

Mentzer *et al.* [124] furthermore trained a PixelCNN-style autoregressive entropy model end-to-end. To soften the non-differentiable discrete rate loss, they used the learned masking technique of Li *et al.* [108], but

formed the surrogate rate loss based on the code length of non-zeroed-out elements of $\hat{\mathbf{z}}$ under the concurrently trained autoregressive entropy model, instead of naive counts as used by Li *et al.* [108].

3.3.3 Uniform Quantization (UQ)

Popularized by Ballé *et al.* [11] and Theis *et al.*, [184], uniform quantization — in its most common form — rounds each element of \mathbf{z} to the closest integer,

$$\llbracket \mathbf{z} \rrbracket := \lfloor \mathbf{z} \rfloor.$$

This can be viewed as a scalar version of the VQ approach, but with a fixed quantization grid equal to the set of integers. The assumption of a uniform quantization grid with width 1 can generally be justified by using a sufficiently flexible pair of transforms (f, g) , which can warp the quantization grid in arbitrary ways if needed [11][10]. Compared to VQ, uniform quantization (and more generally, scalar quantization) is cheap to compute. Moreover, by embedding the integer-valued discrete representation $\hat{\mathbf{z}}$ in \mathbb{R}^N , an entropy model can be conveniently specified in terms of a continuous density model, allowing for simpler differentiable rate surrogates than in approaches based on categorically distributed entropy models (e.g., [1][108][124]). Such an entropy model P is defined by an underlying density p , exactly as in a discretized density model (Eq. 2.5),

$$P(\hat{\mathbf{z}}) := \int_{[-0.5, 0.5]^n} p(\hat{\mathbf{z}} + \mathbf{v}) d\mathbf{v}, \quad \forall \hat{\mathbf{z}} \in \mathbb{Z}^n. \quad (3.12)$$

We defer details on the choice of p to Section 3.3.5, and now discuss a few representative neural compression approaches based on this form of entropy model and integer-valued $\hat{\mathbf{z}}$. In the rest of this sub-section, $\mathbf{u} \sim \mathcal{U}([-0.5, 0.5]^n)$ is a random sample drawn from the uniform density \mathcal{U} on the hypercube $[-0.5, 0.5]^n$.

UQ + STE. Theis *et al.* [184] proposed to train with uniform quantization and approximately differentiate through it by STE, using the identity gradient on the backward pass. For rate control, they optimized the same rate upper bound as on the LHS of Eq. 2.6, replacing the code length $-\log_2 P(\hat{\mathbf{z}})$ by the differentiable upper bound $\mathbb{E}_{\mathbf{u}}[-\log_2 p(\hat{\mathbf{z}} + \mathbf{u})]$.

Additive Uniform Noise. Ballé *et al.* [11] replace rounding with additive uniform noise for model training, i.e.,

$$\lfloor \mathbf{z} \rfloor \approx \mathbf{z} + \mathbf{u}, \quad \mathbf{u} \sim \mathcal{U}([-0.5, 0.5]^n).$$

Naively, one might simply substitute the above into the Lagrangian Eq. 3.10 and hope to obtain a reasonable training objective. However, the resulting code length, $-\log_2 P(f(\mathbf{x}) + \mathbf{u})$, does not yet make sense, as our entropy model P has only been defined over integers. It turns out the form of P (Eq. 3.12) offers a convenient solution: we can simply extend P from \mathbb{Z}^n to all of \mathbb{R}^n , by convolving the underlying density p with the uniform noise \mathbf{u} , i.e.,

$$\tilde{p} := p * \mathcal{U}([-0.5, 0.5]^n). \quad (3.13)$$

It's easy to see that \tilde{p} agrees with P on all integer points, and serves as a smoothed relaxation of P which defines a surrogate gradient with respect to its input. Replacing P by \tilde{p} in Eq. 3.10, and taking expectation with respect to the uniform noise \mathbf{u} , we obtain the surrogate training objective,

$$\mathbb{E}_{\mathbf{x} \sim P_{data}, \mathbf{u} \sim \mathcal{U}}[-\log_2 \tilde{p}(f(\mathbf{x}) + \mathbf{u}) + \lambda \rho(\mathbf{x}, g(f(\mathbf{x}) + \mathbf{u}))], \quad (3.14)$$

which is now differentiable with respect to all components of the model, and can be simply estimated by Monte-Carlo sampling.

Stochastic Gumbel Annealing. Yang *et al.* [209] proposed to optimize the following differentiable surrogate R-D objective:

$$\mathbb{E}_{\mathbf{x} \sim P_{data}} \mathbb{E}_{\hat{\mathbf{z}} \sim q(\hat{\mathbf{z}}|\mathbf{x})}[-\log_2 P(\hat{\mathbf{z}}) + \lambda \rho(\mathbf{x}, g(\hat{\mathbf{z}}))], \quad (3.15)$$

where $q(\hat{\mathbf{z}}|\mathbf{x})$ is an encoding distribution over integer valued latents, to be discussed below. Unlike the usual NELBO, the rate term above is a cross-entropy, rather than a KL divergence. The optimal choice of $q(\hat{\mathbf{z}}|\mathbf{x})$ is therefore deterministic, placing all of its mass on the choice of $\hat{\mathbf{z}}$ that minimizes the R-D cost for each \mathbf{x} , in which case we recover a deterministic R-D VAE (Section 3.2.2). Finding such an optimal deterministic encoder q is a challenging discrete optimization problem, therefore the idea is to relax q into a stochastic encoder (to enable gradient descent), and gradually anneal it towards a deterministic one.

Yang *et al.* [209] parameterized $q(\hat{\mathbf{z}}|\mathbf{x})$ by a continuous location parameter $\mathbf{z} \in \mathbb{R}^n$ (e.g., predicted by an encoder network as $\mathbf{z} = f(\mathbf{x})$) and a temperature hyperparameter $\tau > 0$. It is defined by

$$q(\hat{\mathbf{z}}|\mathbf{x}) := \prod_i q(\hat{\mathbf{z}}_i|\mathbf{x}), \quad (3.16)$$

$$q(\hat{\mathbf{z}}_i|\mathbf{x}) \propto \begin{cases} \exp\{-\psi(\mathbf{z}_i - \lfloor \mathbf{z}_i \rfloor)/\tau\}/C, & \text{if } \hat{\mathbf{z}}_i = \lfloor \mathbf{z}_i \rfloor \\ \exp\{-\psi(\lceil \mathbf{z}_i \rceil - \mathbf{z}_i)/\tau\}/C, & \text{if } \hat{\mathbf{z}}_i = \lceil \mathbf{z}_i \rceil \end{cases} \quad (3.17)$$

where C is a normalizing constant, and $\psi = \tanh^{-1}$. The resulting categorical distribution concentrates its mass on the vertices of a hypercube containing \mathbf{z} , and generalizes stochastic binarization [188] to the integer lattice. The probability mass assigned to each integer neighbor of \mathbf{z} depends inversely on its distance to \mathbf{z} , similar to the softmax formulation of Agustsson *et al.* [1]. The temperature hyperparameter is annealed towards zero over the course of optimization, such that sampling from q converges to deterministic rounding $\mathbf{z} \rightarrow \lfloor \mathbf{z} \rfloor$. For gradient-based optimization, the Gumbel-softmax trick [95][119] is used to differentiate through samples of q , and the discrete entropy model P is replaced by its continuous extension \tilde{p} (Eq. 3.13), as in the uniform noise approach.

This method, Stochastic Gumbel Annealing (SGA), was originally proposed to improve the compression performance of pre-trained models at test time [209]. In this work, \mathbf{z} was initialized to the amortized prediction $f(\mathbf{x})$, but then treated as a variational parameter, and iteratively optimized with gradient descent as in semi-amortized VI [98].

Tsubota *et al.* [196] further applied a version of SGA for end-to-end training, using STE instead of the Gumbel-softmax trick to differentiate through sampling $\hat{\mathbf{z}} \sim q$ (which we refer to as SGA+STE), and obtained improved R-D performance compared to the UQ+STE approach.

Comparisons. Empirical results [106][132][4] suggest that it is beneficial to train with different approximations for optimizing the distortion v.s. the rate terms of the R-D loss (Eq. 3.10). A recent empirical comparison of combinations of various approaches [196] confirms this, showing that it is best to combine a rounding-based approximation (SGA+STE, UQ+STE) for the distortion term, and a uniform-noise-based approximation (additive uniform noise, dithered quantization [36]) for the rate term.

3.3.4 Connection to variational autoencoders, revisited

Previously in Section 3.2.2, we interpreted a general non-linear transform coding model as a rate-distortion VAE (Section 3.1.2) with a discrete latent space and deterministic variational posterior distribution. In this section, we will consider other types of rate-distortion VAEs with continuous latent spaces, and discuss their relation to lossy compression.

First, we show that the additive uniform noise approach for end-to-end training (discussed in Section 3.3.3) equivalently defines a rate-distortion VAE of the data, with a *continuous* latent space and a particular choice of prior and posterior distributions. To derive this, it is instructive to consider the density model p as approximating the distribution of the *continuous* representation $\mathbf{z} = f(\mathbf{x})$ as induced by $\mathbf{x} \sim P_{data}$ and the analysis transform f . Indeed, suppose \mathbf{z} is distributed according to p (if our modeling is perfect), then the distribution of $\hat{\mathbf{z}} = \lfloor \mathbf{z} \rfloor$ has precisely the form of the entropy model P defined by Eq. 3.12. Moreover, if we define the “noisy quantization” by the random variable $\tilde{\mathbf{z}} := \mathbf{z} + \mathbf{u}$, then its induced density (given $\mathbf{z} \sim p$) is precisely \tilde{p} from Eq. 3.13. Based on these connections, the surrogate Lagrangian from additive uniform noise (Eq. 3.14) defines a particular rate-distortion VAE, where $\tilde{\mathbf{z}}$ is the latent variable. Specifically, Eq. 3.14 can be shown [184][12] to be equal to the NELBO objective of an R-D VAE, given by

$$\mathbb{E}_{\mathbf{x} \sim P_{data}} \mathbb{E}_{\tilde{\mathbf{z}} \sim q(\tilde{\mathbf{z}}|\mathbf{x})} [-\log_2 \tilde{p}(\tilde{\mathbf{z}}) + \log_2 q(\tilde{\mathbf{z}}|\mathbf{x}) - \log_2 p(\mathbf{x}|\tilde{\mathbf{z}})] + \text{const}, \quad (3.18)$$

where $\tilde{p}(\tilde{\mathbf{z}})$ plays the role of a prior, $q(\tilde{\mathbf{z}}|\mathbf{x})$ is a fully factorized uniform posterior density centered at $\mathbf{z} = f(\mathbf{x})$, and $p(\mathbf{x}|\tilde{\mathbf{z}})$ is a likelihood density aligned with the distortion function ρ ⁵. The equivalence can be seen by noting that the posterior entropy term is constant (in fact, 0), and sampling from the uniform posterior q is equivalent to adding noise to $f(\mathbf{x})$, by the reparameterization trick.

Although additive uniform noise was originally motivated as a differentiable surrogate to the compression cost using deterministic uniform

⁵In the common case of a squared distortion, $p(\mathbf{x}|\tilde{\mathbf{z}})$ is a Gaussian with mean equal to the reconstruction $\hat{\mathbf{x}} = g(\tilde{\mathbf{z}})$, and covariance inversely proportional to λ . See Section 3.1.2.

quantization (Eq. 3.10), the surrogate loss (Eq. 3.14) can be shown [10][2] to exactly equal the compression cost using *universal quantization*, that is, quantization with a random offset, to be discussed in Section 3.4.1. In other words, given our choice of posterior q and prior \tilde{p} distributions with particular shape restrictions, there is an efficient lossy compression procedure whose communication cost is fully differentiable (in particular, w.r.t. to the encoder parameters), with a bit-rate equal to the (expected) KL divergence $\mathbb{E}_{\mathbf{x} \sim P_{data}}[d_{\text{KL}}[q(\tilde{\mathbf{z}}|\mathbf{x})||\tilde{p}(\tilde{\mathbf{z}})]]$.

Next, we consider the possibility of using other types of R-D VAEs for lossy compression. In generative modeling, the approximate posterior q is typically chosen to be as flexible as possible to achieve a lower NELBO and learn a better model [151]. It is therefore natural to wonder if other choices of the encoding distribution q besides the uniform density can also be used to improve lossy compression.

Given a data point \mathbf{x} and an R-D VAE of the data, one idea is to simply quantize the mode or mean of the variational posterior $q(\mathbf{z}|\mathbf{x})$, as in the uniform quantization approach of Ballé *et al.* [12], and entropy code the resulting quantized representation $\hat{\mathbf{z}}$. However, simple uniform quantization can give poor results, outside of the R-D VAE considered in Eq. 3.18. Yang *et al.* [210] proposed an improved quantization scheme for the case where q is a factorized Gaussian with learned variances across different latent dimensions (as in a standard VAE [99]). In this method, the prior distribution $p(\mathbf{z}_i)$ is used to construct a quantization grid for each latent dimension i , taking inspirations from Arithmetic Coding. This is based on considering k -bit truncations of \mathbf{Z}_i under the probability integral transform, for $k = 1, 2, \dots$, so the resulting set of grid points \mathcal{G}_i consist of nested quantiles (median, quartiles, octiles, etc.) of \mathbf{Z}_i . Then, given a data point \mathbf{x} and its inferred posterior q , the posterior mean μ_i of each dimension i is separately quantized to the corresponding grid \mathcal{G}_i , using a squared error distortion weighted inversely by the posterior variance σ_i (this is essentially $-\log q(\mathbf{z}_i|\mathbf{x}_i)$ for a Gaussian q). The quantized value is found by efficiently solving a discrete optimization problem,

$$\min_{\pi \in \mathcal{G}_i} \frac{(\pi - \mu_i)^2}{\sigma_i^2} + \lambda \mathcal{R}(\pi),$$

where $\mathcal{R}(\pi)$ is the bit-rate associated with grid point π , and λ is a rate-distortion trade-off hyperparameter shared across all latent dimensions. The method therefore assigns more bits to latent dimensions with higher posterior uncertainty/variance, and fewer bits to latent dimensions where $d_{\text{KL}}[q(\mathbf{z}_i|\mathbf{x})||p(\mathbf{z}_i)]$ is small. Furthermore, variable bit-rate compression can be achieved by adjusting λ at compression time, and Yang *et al.* showed that with this approach a single Gaussian VAE [99] can already outperform JPEG in image compression [210]. However, this approach is outperformed by the end-to-end optimized approach [12] at lower bit-rates, and does not operationalize the theoretical rate and distortion losses of the Gaussian VAE.

Another idea is to find ways to transmit a sample of q using close to $\mathbb{E}_{\mathbf{x}\sim P_{data}}[d_{\text{KL}}[q(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]]$ bits (or, close to $d_{\text{KL}}[q(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]$ bits for each given \mathbf{x}). This would allow us to operationalize the R-D loss of the associated R-D VAE, and more generally, yield a lossy compression scheme whose coding cost (as a KL divergence) can often be easily optimized with gradient descent. In theory, this approach can even attain the rate-distortion theoretic limit of lossy compression, up to a logarithmic overhead [211], and can also be more efficient in terms of the rate-distortion-perception trade-off [179], [185]. However, an efficient implementation of this approach is non-trivial, and likely even impossible in the worst case [2]. We will discuss this approach, and the related topics, in Section 3.4.

3.3.5 Entropy models

Various entropy models have been proposed to reduce the bit-rate and improve the rate-distortion performance of neural lossy compression, using largely the same modeling ideas as for lossless compression.

Following quantization, various models can be used to further losslessly compress the resulting discrete representations — for instance, autoregressive models [108][189][124] or off-the-shelf adaptive entropy codes [1]. The most common approach in end-to-end methods is to combine uniform quantization with an entropy model parameterized in terms of a density p as in Eq. 3.12. The simplest choice is a fully-factorized p , resulting in a factorized entropy model. Each marginal of p

is typically parameterized as a mixture distribution [184], or indirectly as the derivative of a deep CDF model [14] (exploiting the relation in Eq. 2.7).

Going beyond factorized entropy models, recent research has explored latent-variable modeling, autoregressive modeling, and their combination, to increase the flexibility of the prior density p and improve the compression bit-rate [14][131][132]. The basic latent variable model approach, commonly referred to as the *hyperprior* approach [14], expresses the entropy model’s underlying density through an additional hierarchy of latent variables \mathbf{h} (“hyper-latents”),

$$\begin{aligned} \mathbf{h} &\sim p(\mathbf{h}), & \hat{\mathbf{h}} &= \lfloor \mathbf{h} \rfloor, \\ \mathbf{z} \mid \hat{\mathbf{h}} &\sim p(\mathbf{z} \mid \hat{\mathbf{h}}), & \hat{\mathbf{z}} &= \lfloor \mathbf{z} \rfloor. \end{aligned}$$

The hyperprior density, $p(\mathbf{h})$, is typically parameterized as in a factorized entropy model, while $p(\mathbf{z} \mid \hat{\mathbf{h}})$ is a density (e.g., factorized Gaussian) whose parameters are predicted from $\hat{\mathbf{h}}$ by a neural network (“hyper-decoder”). Crucially, note that the prior density of \mathbf{z} is conditioned on the discrete $\hat{\mathbf{h}}$, as the hyper-latents must be discretized and entropy-coded first at compression time. The information transmitted in the hyper-latents is known as *side information* [14], and lets the sender and receiver dynamically select an entropy model based on the content of the input data. To train such an entropy model, the rate loss (Eq. 3.9) is modified to account for the side-information, replacing $-\log_2 P(\hat{\mathbf{z}})$ by the joint information content $-\log_2 P(\hat{\mathbf{z}}, \hat{\mathbf{h}}) = -\log_2 P(\hat{\mathbf{h}}) - \log_2 P(\hat{\mathbf{z}} \mid \hat{\mathbf{h}})$; the same techniques from Section 3.3.3 can then be used to differentiate through quantization and rate loss. Empirically, the hyperprior considerably reduces the overall bit-rate of a factorized entropy model, with the side-information comprising a small percentage of the overall rate [14].

The bit-rate can be further improved by additionally modeling $\hat{\mathbf{z}}$ autoregressively similarly to a PixelCNN [131], but results in serial and hence slower decoding. To address this, Minnen *et al.* [132] proposed to instead use channelwise (instead of spatial) autoregressive conditioning, significantly speeding up (de)compression without harming the rate-distortion performance. Compared to an autoregressive model, a latent-variable entropy model has the advantage of parallel encoding/decoding

via (conditionally) factorized distributions, but entails transmitting side-information, similar to the two-part code in lossless compression (Section 2.3.1). Yang *et al.* [209] applied bits-back coding to reduce the transmission of side-information in a hyperprior model.

Regardless of the choice of an entropy model, the sender and receiver must agree on the exact same probabilities for entropy coding (such as Arithmetic Coding, discussed in Section 2.1.4) to operate correctly. This can be a stringent requirement when the entropy models (e.g., the conditional model $P(\hat{\mathbf{z}}|\hat{\mathbf{h}})$) are computed on the fly, especially in the face of round-off errors from floating point arithmetic and non-deterministic GPU operations. We refer readers to Ballé *et al.* [13] for more details on this issue, and a potential solution based on integer arithmetic.

3.4 Compression without quantization

The non-differentiability of quantization has hindered end-to-end training of lossy compression models. The various methods in Section 3.3 replace quantization by a differentiable surrogate at training time, leading to a mismatch between “soft” quantization during training and “hard” quantization at test time. This mismatch generally results in sub-optimal performance [209]. Annealing can alleviate the problem [1][209][2] but may suffer high-variance gradients as the approximation approaches hard quantization, and requires specifying an annealing schedule.

Perhaps surprisingly, quantization can be avoided entirely if we are willing to accept some noise in the transmitted representation. Instead of quantization, it is possible to transmit a continuous but stochastic sample $\mathbf{z} \sim q(\mathbf{z} | \mathbf{x})$ using a finite number of bits [18]. The problem of efficiently communicating such a sample is also recognized as *channel simulation* [215][107], *reverse channel coding* [18][2][186], or *relative entropy coding* [57][58].

Bits-back coding (Section 2.3.2) seems like a natural candidate for this problem as it also uses a stochastic encoder. Unfortunately, a requirement of bits-back coding is that the exact data \mathbf{x} is eventually available to the decoder and is therefore only directly applicable to lossless compression. That is, bits-back is a solution to the lossless source coding problem but not a solution for reverse channel coding.

Li *et al.* [107] showed that it is possible to communicate \mathbf{z} at an average coding cost of at most

$$I[\mathbf{X}, \mathbf{Z}] + \log_2(I[\mathbf{X}, \mathbf{Z}] + 1) + 5$$

bits. That is, the coding cost is close to the information contained in \mathbf{z} . However, they also showed that in general it is not possible to significantly reduce this coding cost further. Even for optimal encoders and decoders we may therefore have to pay an overhead which is logarithmic in the mutual information. Note that this overhead is relatively small if the mutual information is large.

One way to increase the mutual information (and thus reduce the relative overhead) is to communicate more information at once (for example, by bundling multiple frames of a video). Unfortunately, it can be computationally very expensive to do so. Agustsson & Theis [2] showed that there is no general reverse channel coding algorithm whose computational cost is polynomial in the information content. If we want to transmit large amounts of information at once using as few bits as possible, then this may only be possible by spending a lot of computation. Nevertheless, some distributions can be communicated efficiently, both computationally and with low overhead.

In the following, we will review two strategies for communicating stochastic information. One is a simple and efficient approach for simulating channels with additive uniform noise, and one is a general approach for communicating samples of arbitrary distributions. For a more thorough introduction to reverse channel coding, see Theis & Yosri [186].

3.4.1 Dithered quantization

Consider a latent representation \mathbf{Z} ⁶ which is the output of a neural network followed by additive uniform noise, i.e.,

⁶To reduce clutter, our notations here differs from those in Section 3.2 on non-linear transform coding. Here, we denote the noise-injected latent representation by \mathbf{z} (rather than $\tilde{\mathbf{z}}$), and denote transform coefficients by \mathbf{y} (rather than \mathbf{z}). Our use of \mathbf{z} here is consistent with Section 3.1.2 and broader machine learning literature on latent variable models (see [99], [217]).

$$\mathbf{U} \sim \mathcal{U}([-0.5, 0.5]^n), \quad (3.19)$$

$$\mathbf{Y} = f(\mathbf{X}), \quad (3.20)$$

$$\mathbf{Z} = \mathbf{Y} + \mathbf{U}. \quad (3.21)$$

It turns out that we can efficiently communicate an instance of \mathbf{z} using an old technique called *dithered* or *universal quantization* [155][223].

Let \mathbf{U}' be another vector of uniform noise independent of \mathbf{U} which is available to both the encoder and the decoder. In practice this requirement can be achieved by generating noise using a pseudorandom number generator with the same random seed. For any fixed value of $\mathbf{y} = f(\mathbf{x})$, it holds that [223],

$$\lfloor \mathbf{y} - \mathbf{U}' \rfloor + \mathbf{U}' \sim \mathbf{y} + \mathbf{U}. \quad (3.22)$$

That is, subtracting uniform noise, rounding, and then adding uniform noise back is distributionally equivalent to adding noise directly. We can exploit this for the communication of a uniform sample as follows. Define the random variable $\mathbf{K} = \lfloor \mathbf{Y} - \mathbf{U}' \rfloor$. Given a data sample \mathbf{x} , the encoder computes $\mathbf{y} = f(\mathbf{x})$, samples a value of \mathbf{u}' , computes $\mathbf{k} = \lfloor \mathbf{y} - \mathbf{u}' \rfloor$, and entropy encodes \mathbf{k} into bits. The decoder receives \mathbf{k} and simply adds \mathbf{u}' , and effectively obtains a sample of $\mathbf{y} + \mathbf{U}$ (by Eq. 3.22). To entropy encode \mathbf{k} , we need to know its distribution. Assuming \mathbf{Z} has marginal distribution $p_{\mathbf{Z}}$, we have [216]

$$P(\mathbf{K} = \mathbf{k} \mid \mathbf{U}' = \mathbf{u}') = p_{\mathbf{Z}}(\mathbf{k} + \mathbf{u}'). \quad (3.23)$$

Note that we can condition on \mathbf{u}' to encode \mathbf{k} since it is known to the decoder. Eq. 3.23 tells us that in order to encode \mathbf{k} , we only need a model for the density of \mathbf{Z} . The expected coding cost is [2], [216]

$$\mathbb{E}[-\log p_{\mathbf{Z}}(\mathbf{K} + \mathbf{U}')] = \mathbb{E}[-\log p_{\mathbf{Z}}(\mathbf{Y} + \mathbf{U})] \quad (3.24)$$

$$\geq h[\mathbf{Z}] \quad (3.25)$$

$$= h[\mathbf{Z}] - \underbrace{h[\mathbf{Z} \mid \mathbf{X}]}_{=0} \quad (3.26)$$

$$= I[\mathbf{X}, \mathbf{Z}] \quad (3.27)$$

with equality when $p_{\mathbf{Z}}$ is the true marginal distribution of \mathbf{Z} .

This coding cost has two useful properties. First, it is equal to the amount of information transmitted (assuming $p_{\mathbf{Z}}$ models \mathbf{Z} faithfully).

That is, encoding \mathbf{K} is a statistically efficient strategy for communicating a sample. Second, Eq. 3.24 is differentiable in \mathbf{y} so that we can easily optimize an encoder using backpropagation. Agustsson & Theis [2] exploited these facts to train neural encoders and decoders for images without the train-test mismatch commonly introduced by quantization.

3.4.2 Minimal random coding

While dithered quantization can be computationally and statistically efficient, it is only able to communicate certain simple distributions. Several general algorithms have been developed to communicate a sample from arbitrary distributions [75][41][42][107][76], though some have only been studied in the context of discrete distributions. Here we describe one algorithm based on importance sampling. In information theory, it is known as the *likelihood encoder* [42][171]. In machine learning, it has recently been introduced as *minimal random coding* (MRC) [76][57].

Assume both the encoder and decoder have access to $p(\mathbf{z})$. For each given data observation \mathbf{x} , the encoder generates N examples $\mathbf{z}_n \sim p(\mathbf{z})$ and forms importance weights for a target distribution $q(\mathbf{z} | \mathbf{x})$,

$$w_n = \frac{q(\mathbf{z}_n | \mathbf{x})}{p(\mathbf{z}_n)}. \quad (3.28)$$

It then randomly samples an index k using the normalized importance weights, i.e.,

$$P(n) = \frac{w_n}{\sum_{i=1}^N w_i}, \quad k \sim P. \quad (3.29)$$

The index is uniformly distributed and encoded using $\log_2 N$ bits. The decoder receives k and reconstructs \mathbf{z}_k . It can do this if, for example, the encoder used a pseudorandom number generator to generate \mathbf{z}_n with a random seed known to the decoder. The complexity of this procedure thus depends on the number of samples N that needs to be generated.

Havasi *et al.* [76] showed that if the number of samples is

$$N = 2^{d_{\text{KL}}[q(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z})] + t} \quad (3.30)$$

and if the distribution of $\log q(\mathbf{Z} | \mathbf{x})/p(\mathbf{Z})$ is concentrated around its expected value, then \mathbf{Z}_k quickly converges to \mathbf{Z} (in a total variation

sense) as t increases. Havasi *et al.* [76] applied MRC to neural network compression, while Flamich *et al.* [57] used it for image compression. Theis & Yosri [186] showed that the coding cost of minimal random coding can be further reduced without any loss in quality.

Cuff [41] considered the setting where m data points are encoded and communicated at once, and $p(\mathbf{z})$ is the marginal distribution of \mathbf{Z} , that is, the average of the distributions $q(\mathbf{z} | \mathbf{x})$ if we average over all $\mathbf{x} \sim p_{data}(\mathbf{x})$. In this setting, he showed that if

$$N > 2^{mI[\mathbf{X}, \mathbf{Z}]},$$

then the distribution of $(\mathbf{Z}_k, \mathbf{X})$ converges to $q(\mathbf{z} | \mathbf{x})p_{data}(\mathbf{x})$ in total variation distance as m goes to infinity. In other words, using on average $I[\mathbf{X}, \mathbf{Z}]$ bits, we can communicate a sample which approximately follows $q(\mathbf{z} | \mathbf{x})$ in a total variation sense.

3.4.3 Stochastic versus deterministic coding

Communication of information without any quantization requires a certain level of noise to be present. Without noise or quantization, there would be no limit to the amount of information we could send through the bottleneck of an autoencoder. This raises the question of whether a deterministic encoder with quantization or a stochastic encoder is better. Ballé *et al.* [10] argue that we can always improve on dithered quantization with a deterministic encoder when performance is measured by a rate-distortion trade-off. Theis & Agustsson [179] extended this argument to arbitrary stochastic encoders. That is, when we care about a rate-distortion trade-off, the best stochastic encoder is likely to perform worse than the best deterministic encoder. However, Theis & Agustsson [179] also showed by example that when we additionally care about the realism of reconstructions (discussed in Section 3.5), stochastic encoders can perform significantly better than deterministic ones. Which one is better therefore depends on the setting of interest.

3.4.4 Lossy compression with diffusion

Reverse channel coding schemes enable novel compression schemes which deviate from the classical transform coding setup. One such approach

is based on Gaussian diffusion models [168]. These generative models learn the joint distribution of the data $\mathbf{x} = \mathbf{z}_0$ and copies of the data corrupted by increasing levels of noise,

$$\mathbf{z}_{t+1} = \alpha \mathbf{z}_t + \beta \mathbf{v}_t, \quad (3.31)$$

where \mathbf{v}_t is Gaussian noise. For appropriate α and β , \mathbf{z}_t as a function of t can be viewed as approximating a Gaussian diffusion process.

Ho *et al.* [82] first considered the rate-distortion performance of a scheme which would efficiently communicate an instance of \mathbf{z}_t for some fixed t before generating an estimate of $\mathbf{x} = \mathbf{z}_0$ with the help of the diffusion model. Theis *et al.* [185] demonstrated that this approach can work extremely well on small images when compared to the best transform coding schemes, especially when realism is considered. This is surprising considering that Gaussian noise is added directly to the data instead of applying an encoder transform first. Unlike neural transform coding which typically requires training many different models, the same diffusion model can be used to encode and decode a Gaussian sample at arbitrary bit-rates. The simplicity of this approach makes it very attractive from a theoretical perspective but its practicality is still an open question as the approach is very computationally expensive.

Diffusion models have also been used in a transform coding setting by conditioning the generative model on the quantized output of an encoder. Saharia *et al.* employed a diffusion model for artefact removal on JPEG images [159], while Yang *et al.* [207] conditioned a diffusion model on a discrete latent variable in an end-to-end trained variational autoencoder framework [207]. These approaches do not require the as of yet computationally expensive communication of a sample but require dealing with a non-differentiable quantization operation as in other transform coding schemes. They also do not benefit from the potential bit-rate savings of coding with a stochastic encoder [179][185] (Section 3.4.3).

3.5 Perceptual losses

Neural networks are only as good as the losses they are trained for. While in lossless compression the objective is clear – namely to minimize

the required number of bits to represent the data – the story is a lot more complicated when we turn to lossy compression. Here we need to make decisions about which information to sacrifice in order to save additional bits. In typical media compression applications, the goal is to make any reconstruction errors as imperceptible as possible, which raises complicated questions about how our brains perceive differences between signals.

3.5.1 Background

We can distinguish between two types of distortions, namely *full-reference metrics* and *no-reference metrics*. The former is a function which takes an image and its reconstruction as input while the latter only looks at the reconstruction to make a judgement about its quality. These can be motivated by two corresponding types of quality measures involving humans. *Mean-opinion scores* (MOS) [93] are measured by having raters judge reconstructions on a scale from 1 (bad) to 5 (excellent). Since raters are only provided with the reconstructions, their judgments can be viewed as the output of a no-reference metric. Many other *image quality assessments* (IQAs) evaluate perceptual quality without reference to the original data [174], but MOS is the most widely used measure due to its simplicity. In contrast, *degradation MOS* (DMOS) asks raters for a judgment based on both the unprocessed and the reconstructed data [93], and thus is more akin to a full-reference metric.

In addition to distortions, we may consider *divergences* which depend on the probability distribution of the data and the marginal distribution of reconstructions. Driven mostly by the success of *generative adversarial networks* (GANs) [65] in producing realistic looking images, divergence optimization has become an important topic in neural compression.

3.5.2 Perceptual distortions

In this section we explore some of the metrics that are most likely to be encountered in the current literature on neural compression. However, there is a much larger body of potentially relevant work on IQAs and more broadly on perception that we will have to ignore here. For

example, *VMAF* [113][109] is a full-reference metric which is frequently used to evaluate the quality of video in industry but it has not yet found widespread use in the neural compression community.

Mean-squared error (MSE) and peak signal-to-noise ratio (PSNR) are frequently used in neural compression but do not predict perceived quality well [55]. Another criterion widely used to evaluate images in machine learning and beyond is the structural similarity index (SSIM) [221]. SSIM has been shown to correlate better with human visual perception and has been extensively studied and extended. The extension most commonly found in neural compression papers is the multi-scale SSIM (MSSSIM) [201] which evaluates SSIM at multiple resolutions and combines them multiplicatively. The neurogram similarity index measure (NSIM) is another closely related metric which is applied to neurograms or spectrograms of audio signals [79] and has been shown to correlate well with the perceived quality of speech [80].

While MSE and PSNR are computed pixel-wise, SSIM is computed from small image patches. Let \mathbf{x} and \mathbf{y} be two aligned grayscale image patches extracted from an image and its reconstruction, respectively. Further, let $\mu_{\mathbf{x}}, \sigma_{\mathbf{x}}$, and $\sigma_{\mathbf{xy}}$ represent the average pixel value, the standard deviation and the covariance of pixel values as measured from these patches. Using these quantities, we define

$$l(\mathbf{x}, \mathbf{y}) = \frac{2\mu_{\mathbf{x}}\mu_{\mathbf{y}} + C_1}{\mu_{\mathbf{x}}^2 + \mu_{\mathbf{y}}^2 + C_1}, \quad (3.32)$$

$$c(\mathbf{x}, \mathbf{y}) = \frac{2\sigma_{\mathbf{x}}\sigma_{\mathbf{y}} + C_2}{\sigma_{\mathbf{x}}^2 + \sigma_{\mathbf{y}}^2 + C_2}, \quad (3.33)$$

$$s(\mathbf{x}, \mathbf{y}) = \frac{\sigma_{\mathbf{xy}} + C_3}{\sigma_{\mathbf{x}}\sigma_{\mathbf{y}} + C_3}, \quad (3.34)$$

where C_1, C_2, C_3 are positive constants to ensure numerical stability. The structural similarity function s measures correlation, while the luminance function l and contrast function c are chosen to respond to *relative* changes in luminance or contrast. These functions will not change much if, for example, both $\mu_{\mathbf{x}}$ and $\mu_{\mathbf{y}}$ are scaled by the same factor. This is consistent with Weber's law of how the human visual system perceives changes in these parameters [201]. SSIM is defined as

$$\text{SSIM}(\mathbf{x}, \mathbf{y}) = l(\mathbf{x}, \mathbf{y})^\alpha c(\mathbf{x}, \mathbf{y})^\beta s(\mathbf{x}, \mathbf{y})^\gamma, \quad (3.35)$$

where α, β, γ are additional parameters which control the relative importance of the factors (typically set to 1). Note that $\text{SSIM}(\mathbf{x}, \mathbf{x}) = 1$. To compute a single value for an entire image, one approach is to use a sliding window (e.g., 8×8 pixels) and then to average SSIM values. MS-SSIM instead uses a smooth windowing approach to compute local statistics in order to avoid blocking artifacts [201]. SSIM is defined for grayscale images. In the neural compression literature, SSIM is typically applied separately to RGB channels and the resulting values are averaged to evaluate color images.

SSIM has been shown to perform significantly better at predicting human judgments than MSE on common distortions such as blurriness, noise or blocking artifacts [221]. However, its limitations are also well documented and it tends to fail for reconstructions produced by generative compression approaches [105][127].

“When a measure becomes a target, it ceases to be a good measure” [175]. In line with this adage, directly optimizing neural networks for MS-SSIM offers mixed results when compared to MSE in terms of perceptual quality [14]. Unlike many applications of IQA, a metric has to make meaningful predictions for all conceivable distortions to be useful as a target in neural compression and cannot have any blind spots. Ding *et al.* [45] recently compared a large number of IQA methods and found that many were unsuitable for direct optimization.

A common approach to the design of more sophisticated distortions is to rely on deep neural networks which perform well in some other vision task. Typically, these distortions take the form

$$\rho_{\Phi}(\mathbf{x}, \mathbf{y}) = \rho(\Phi(\mathbf{x}), \Phi(\mathbf{y})), \quad (3.36)$$

where Φ is some representation derived from the hidden activations of a neural network and ρ is typically the MSE. Gatys *et al.* [62] compellingly demonstrated the ability of such metrics to capture semantic content at different levels of abstraction in their seminal paper on neural style transfer. Bruna *et al.* [30] used distortions derived from VGG [166] and scattering networks [31] and applied them to the task of super-resolution, which can be viewed as a simpler form of neural compression with a fixed encoder. They found that these distortions lead to sharper reconstructions than MSE but can also cause artifacts.

Zhang *et al.* [219] further investigated the efficacy of distortions based on VGG and found that they can significantly outperform SSIM on a range of artifacts including those generated by neural networks. They further proposed the *learned perceptual image patch similarity* (LPIPS). LPIPS uses pretrained classifiers such as AlexNet [103] or VGG [166] but its parameters are finetuned in a supervised manner to match human responses. As an alternative, Bhardwaj *et al.* [21] recently showed that representations learned in a completely unsupervised manner can be as effective as LPIPS and proposed the *perceptual information metric* (PIM). Here, Φ was learned using a contrastive loss.

As of this writing, no metric has reached the level of humans predicting the responses of other humans and how to close the gap is an open research question. Amir & Weiss [8] found that randomly initialized networks and a simple kernel-based metric can perform as well as VGG-based distortions in predicting human responses, raising interesting questions about the necessity and usefulness of neural representations in perceptual distortions.

3.5.3 No-reference metrics

While no-reference metrics are rarely used as training targets in the current neural compression literature, they are sometimes used for evaluation [127]. The *natural image quality evaluator* (NIQE) [134], for example, extracts nonlinear features from image patches sampled from a test image. It then fits a Gaussian distribution to these features and compares it to a Gaussian distribution fitted to features extracted from natural images. Many extensions of NIQE have been proposed and used in the IQA literature [218].

Deep IQA [28] samples 32x32 image patches from a test image and uses a convolutional neural network to predict perceptual quality judgments. The predictions for different patches are averaged to yield a single score for an image. The parameters of the network were trained in a supervised manner on the LIVE dataset [165] which only contains 981 distorted versions of 29 reference images. To augment this dataset and to reduce overfitting, Kim *et al.* [97] pretrained a neural network to predict pixel-wise distortions before training on the LIVE dataset.

3.5.4 Divergences and adversarial networks

An image compressor which always outputs a fixed image of high perceptual quality would score high in terms of any no-reference metric and consume zero bits. However, for a useful compressor of natural images, we expect a diverse set of reconstructions roughly following the distribution of uncompressed images. Such properties can be assessed with *divergences*, measuring the deviation between the data distribution and the distribution of reconstructions. Requirements on a divergence d are $d[p, q] \geq 0$ and $d[p, q] = 0 \iff p = q$ for all distributions p, q . A divergence does not need not be symmetric. Examples include the Kullback-Leibler divergence (KLD), the Jensen-Shannon divergence (JSD), or the total variation distance (TVD).

Divergences are sometimes described as either *zero-avoiding* or *zero-forcing* [130]. Zero-avoiding divergences assign a high penalty to models which assign zero probability to events ($q(\mathbf{x}) = 0$) which have positive probability under the reference distribution ($p(\mathbf{x}) > 0$). They are also called *mode-covering* divergences since they encourage models to assign probability mass to all modes of a distribution. Examples include the KLD or the χ^2 -divergence. On the other hand, zero-forcing divergences assign a high penalty to model distributions which assign positive probability ($q(\mathbf{x}) > 0$) to events which have zero probability ($p(\mathbf{x}) = 0$). These are also called *mode-seeking* divergences since the resulting models tend to ignore some of the modes of distributions with multiple peaks. An examples is the reverse KLD (Eq. 2.4 with P and Q switched). Zero-forcing divergences are especially useful for capturing realism since they discourage models from generating implausible reconstructions.

We can give further motivation for the total variation distance (TVD). Instead of asking human observers to rate reconstructions as in a MOS test, consider an experiment where we randomly present either real data or reconstructions with equal probability, and ask the observers to make a binary decision as to whether the data shown is real [44]. An optimal observer will correctly classify the data with probability [25][137]

$$p_{\text{success}} = \frac{1}{2}d_{\text{TV}}[q, p] + \frac{1}{2}, \quad (3.37)$$

where d_{TVD} is the TVD,

$$d_{\text{TVD}}[q, p] = \int |q(\mathbf{x}) - p(\mathbf{x})| d\mathbf{x}. \quad (3.38)$$

That is, minimizing TVD minimizes the chance of an optimal classifier correctly discriminating reconstructions from real data. Other divergences can be motivated by considering other losses and the associated risk of an optimal classifier [137][172]. Optimizing divergences of high dimensional distributions is challenging. Nevertheless, approximations optimized by generative adversarial networks (GANs) have proven useful in practice [65][138].

Weighted combinations of distortions and adversarial losses have produced very promising results in neural compression and related image reconstruction tasks [105][152][161][3][127],

$$\sup_{D \in \mathcal{D}} \mathcal{L}_D(f, g) + \beta \mathcal{L}_\rho(f, g), \quad (3.39)$$

where D is a discriminator network (adversary), and (f, g) are a pair of encoder and decoder networks introduced in Section 3.2.1. To give a simple example, we might choose

$$\mathcal{L}_D(f, g) := \mathbb{E}[\log D(\mathbf{X}) + \log(1 - D(g(\llbracket f(\mathbf{X}) \rrbracket), \epsilon))], \quad (3.40)$$

$$\mathcal{L}_\rho(f, g) := \mathbb{E}[|\Phi(\mathbf{X}) - \Phi(g(\llbracket f(\mathbf{X}) \rrbracket), \epsilon)|^2]. \quad (3.41)$$

The feature representation Φ is typically a combination of pixels and VGG feature activations. $\llbracket \cdot \rrbracket$ is a quantizer mapping the output of the encoder to a discrete number of values. In addition to the encoder's outputs, the decoder receives independent noise ϵ as input. It can be shown that $\mathcal{L}_D(f, g)$ lower-bounds the Jensen-Shannon divergence between the data distribution and the marginal distribution of reconstructions [65]. Similar losses can be used to target other divergences [138]. For example, Agustsson *et al.* [3] used an LSGAN loss which targets a χ^2 -divergence [121] and were able to train autoencoders achieving much more detailed reconstructions at extremely low bit-rates than advanced classical codecs.

Divergences are not just used for training but also for evaluating neural compression results. In particular, the Frechet inception distance (FID) [78] measures the squared Wasserstein-2 distance between the

data and reconstructions in the feature space of the *inception-v3* network [177], approximating both distributions as multivariate Gaussian. Both FID and NIQE (Section 3.5.3) measure the distance between two Gaussian distributions. However, while NIQE estimates the distribution of features *within* a single image, FID measures the distribution of features *across* many different images.

3.5.5 Perception-distortion trade-off

An important question for neural compression is whether divergences are needed at all. Could we achieve the same results with just a distortion? Optimizing MSE, SSIM and even neural-network-based losses tend to produce artefacts but perhaps these will disappear as perceptual distortions improve. A strong counterargument to this hope was given by Blau & Michaeli [25]. They showed that *any* distortion is going to produce artefacts for some data unless the compressed representation allows us to reproduce the inputs perfectly. That is, for any lossy codec, the optimal reconstruction $\hat{\mathbf{X}}$ of \mathbf{X} optimized for the expectation of a given distortion will not preserve the data distribution [25, Theorem 1]. Note that this limitation does not apply if we are allowed to optimize divergences since then we can achieve perfect realism by requiring that all admissible decoders produce reconstructions such that $d[p_{\hat{\mathbf{X}}}, p_{\mathbf{X}}] = 0$ for some divergence d , at least in theory. Blau & Michaeli [25] also showed that the achievable divergence can only increase when we demand lower levels of distortion. This makes sense, as the set of available decoders shrinks as we impose stronger constraints on them. Counterintuitively, however, this means that minimizing distortion can have the effect of reducing the realism of reconstructions. Blau & Michaeli [25] called this the *perception-distortion trade-off*, that is, a small distortion limits the achievable realism and a low divergence (high realism) limits the achievable distortion. However, the amount of tension between distortions and divergences depends on their particular choices [50] [105] [25].

3.6 Task-oriented compression

An increasing amount of data, such as climate data or satellite images, require efficient storage and transmission, and will only ever be “seen” by algorithms and machines that process them [53]. For such data, the “perceptual” quality of the reconstructions is irrelevant; rather, we only wish to preserve our performance on certain downstream tasks, when we use a compressed representation instead of the original data. By performing downstream tasks directly on the compressed representation, and instead of a reconstruction in the data space, we can potentially achieve savings both in the bit-rate [167][53] as well as computational requirements of the system [190][123].

We may formalize *task-oriented compression* as follows. Suppose we are provided (\mathbf{X}, \mathbf{Y}) pairs from some joint distribution, where \mathbf{X} needs to be compressed and sent to the receiver as before, while \mathbf{Y} is available to both the sender and receiver, and has the interpretation of a *target* or supervision signal for some task. A concrete example could be semantic segmentation, where \mathbf{x} is an image, and \mathbf{y} contains ground-truth semantic category labels for each pixel of \mathbf{x} . Given a sampled pair of data and target (\mathbf{x}, \mathbf{y}) , a task-oriented compression algorithm maps \mathbf{x} to a representation \mathbf{z} , either deterministically or stochastically (in either case, \mathbf{z} is another random variable constructed from \mathbf{x}). The transmission of \mathbf{z} incurs some bit-rate loss \mathcal{R} , as usual. The decoder, having received \mathbf{z} , performs some downstream computation and evaluates performance given the target \mathbf{y} , resulting in a task-oriented distortion loss \mathcal{D} . Task-oriented compression then aims to minimize a combination of the rate and task-oriented distortion losses, with different methods making different choices of \mathcal{R} and \mathcal{D} .

Arguably the earliest example is the Information Bottleneck principle [187]. Here the rate loss is given by the mutual information $\mathcal{R} = I[\mathbf{X}, \mathbf{Z}]$ as in the rate-distortion function (Section 3.1), while the distortion is chosen to be a negative mutual information $\mathcal{D} = -I[\mathbf{Z}, \mathbf{Y}]$. The general idea is to encourage \mathbf{Z} to become a minimal sufficient statistic of \mathbf{X} for predicting \mathbf{Y} [7].

Practical examples of task-oriented compression usually require \mathbf{Z} to be discrete, so it can be compressed with the neural lossy compression

approaches described in Section 3.2. The rate loss is thus measured by the Shannon entropy, $\mathcal{R} = H[\mathbf{Z}]$, and the distortion is based on a concrete task, instead of an information-theoretic one as in the Information Bottleneck. Often, the decoder further computes a prediction $\hat{\mathbf{y}}$ given \mathbf{z} , e.g., by passing \mathbf{z} through a neural network. The task distortion may be a conditional cross-entropy, $\mathcal{D} = H[\mathbf{Y}, \hat{\mathbf{Y}}|\mathbf{Z}]$, for a classification task⁷, or a squared error loss, $\mathcal{D} = \mathbb{E}[\|\mathbf{Y} - \hat{\mathbf{Y}}\|^2] = \mathbb{E}[\|\mathbf{Y} - g(\mathbf{Z})\|^2]$, for a regression task. Singh *et al.* [167] considered the classification setting, where \mathbf{Z} is a (often high-dimensional) feature tensor produced by the penultimate layer of a deep CNN. They optimized for the bit-rate and classification accuracy by training end-to-end on the Lagrangian objective $\mathcal{R} + \lambda\mathcal{D}$, using the same additive uniform noise approximation as in [12] (discussed in Section 3.3.3). Matsubara *et al.* [123] adopted a similar end-to-end approach, but \mathbf{Z} is chosen to be computed by an earlier layer of a neural network, due to a limited computational budget of the sender (e.g., an edge device sending a photo to a cloud server). Dubois *et al.* [53] extended the task-oriented compression setup to consider multiple downstream tasks, aiming to learn a compressed representation that ensures good performance on a variety of tasks. A challenge is that we rarely know in advance all the downstream tasks of future interest, at compression time. They address this by focusing on tasks that are invariant under user-defined transformations to the input data \mathbf{x} , such as image classification under random translation or cropping.

3.7 Video compression

While deep generative modeling has impacted image compression early on [70], its application to *video* compression began more recently (circa 2018) due to the additional complexity in modeling videos as well as the increased computational complexity.

⁷Let $p_{\mathbf{ZY}}$ denote the joint distribution of \mathbf{Z} and \mathbf{Y} , and suppose the decoder makes a probabilistic classification $\hat{\mathbf{y}} \sim q_{\hat{\mathbf{Y}}|\mathbf{Z}=\mathbf{z}}$ for any given \mathbf{z} . The task loss for maximizing the likelihood of \mathbf{Y} given \mathbf{Z} is then $\mathcal{D} = \mathbb{E}_{(\mathbf{z}, \mathbf{y}) \sim p_{\mathbf{ZY}}}[-\log q_{\hat{\mathbf{Y}}|\mathbf{Z}}(\mathbf{y}|\mathbf{z})] = \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{Z}}}[H[p_{\mathbf{Y}|\mathbf{Z}=\mathbf{z}}, q_{\hat{\mathbf{Y}}|\mathbf{Z}=\mathbf{z}}]] =: H[\mathbf{Y}, \hat{\mathbf{Y}}|\mathbf{Z}]$, where $H[\cdot, \cdot]$ denotes (unconditional) cross entropy (defined in Eq. 2.3).

A typical video codec consists of two steps: motion compensation and residual compression. The idea behind motion compensation is to predict the next frame in a video based on previous frames. Traditionally, motion compensation relied on block motion estimation (i.e., matching entire patches in videos and memorizing compressed displacement vectors). By contrast, neural approaches typically directly predict the displacement fields of pixels between adjacent frames. If the predicted displacement field is “simple”, e.g., sparse, it can be efficiently compressed. Since certain aspects in a video can not easily be predicted, the (typically sparse) residual is separately compressed using image compression models (classical or neural).

Recently proposed neural video codecs differ in design choices of the predictive model (e.g., stochastic vs deterministic) and the residual compression scheme (e.g., compression in latent space vs. pixel space). Another fundamental design choice is to either consider the low-latency setup in which a video has to be encoded and decoded on a frame-by-frame level, or the offline setup, in which case the video is encoded as a whole (using knowledge of future frames). For example, the offline setup may be adequate for video streaming services, while the online setup may be more suitable for video conferencing. Another line of research investigates hybridizing classical and neural codecs [117].

One of the earliest approaches to joint prediction and residual compression based on convolutional architectures was [33], which still used a traditional block-based motion estimation approach. For offline compression, [203][48] formulated the video compression problem as a frame interpolation problem. Here, a subset of “key frames” are compressed as images, and the intermediate frames reconstructed based on video interpolation techniques.

Most neural compression approaches currently focus on the low-latency (online) setup [74][115][4][64][206][116]. These approaches can be interpreted as autoregressive generative models for frame sequences [208]. Han *et al.* [74] and Habibian *et al.* [73] first adapted the neural image compression framework of Ballé *et al.* [12] to video data, framing it in a latent variable modeling context. While Han *et al.* [74] proposed to encode the frame sequence using the predictive next-frame distribution of a stochastic recurrent neural network (eliminating the need

to separately compress residuals) and adding a global latent variable that played a similar role as a key frame in traditional compression, Habibian *et al.* [73] used 3D convolutions and explored structured priors such as PixelCNNs [141] for lower bit-rates at the expense of increased runtime.

In subsequent work, the separation of motion estimation and residual compression dominated and led to improvements over classical video codecs such as HEVC. Lu *et al.* [115] adopted a hybrid architecture that combined a pre-trained Flownet [51] and residual compression. Rippel *et al.* [153] proposed a motion estimation approach with long-term memory and adaptive rate control.

Another noteworthy innovation is scale-space flows. Agustsson *et al.* [4] used learned optical flows and warping to predict the next frame in sequence, however, it does so by adding a “scale” dimension to the optical flow field. This dimension allows the model to adaptively blur the source based on how well the next frame is predicted. A general framework for low-latency video compression was recently introduced by Yang *et al.* [208]. The paper also showed that the frame reconstruction for models such as those of Agustsson *et al.* [4] could be improved by introducing a scale parameter that mediates between the autoregressive prediction and the compressed residual (in a similar way as how RealNVP improves over NICE [46][47]), as well as by using non-factorized priors.

The components of neural video compression models are often strongly inspired by classical codecs, in particular those dealing with motion compensation. While this allows neural models to approximate the performance of their classical counterparts, it also introduces complexity and can limit their flexibility. Recent work was able to achieve state-of-the-art results using a greatly simplified approach dubbed the video compression transformer (VCT) [126]. Here, individual frames are independently transformed by encoder and decoder transforms but the latent representations are jointly encoded with a powerful entropy model based on transformers [198].

4

Discussion and Open Problems

Neural compression is a rapidly growing field and has made significant strides in both lossless and lossy compression. While neural approaches to image compression barely beat JPEG 2000 in 2016 [184][12], they already outperformed the best known handcrafted codecs in 2018 [220]. And in the *Challenge on Learned Image Compression*¹ of 2021, classical codecs did not even make it into the top 10. Similarly, the leading codec in the *Large Text Compression Benchmark* relied solely on neural networks to predict text [16].

Nevertheless, many practical and theoretical challenges remain. Chief among them is computational complexity, which stands in the way of the wider adoption of neural compression. Computational feasibility is also a critical factor in the application of neural networks to new data modalities, such as point clouds and VR content. There are also many open problems in neural lossy compression in particular, such as the design of loss functions and evaluation criteria, more efficient compression without quantization, and ways to mitigate the risk of miscommunication.

¹<http://compression.cc/>

A major roadblock for both neural lossless and lossy compression is high computational complexity. While neural-network-based approaches offer remarkable compression performance, they demand significantly higher computation than traditional codecs, and have so far mostly been developed in high-performance computing environments with GPUs. However, real-world applications, such as video streaming on mobile devices, come with stringent requirements on latency, power consumption, hardware compatibility, etc., at a much lower computation budget. Neural compression methods will need to meet these requirements, and still deliver significant improvements over traditional methods, for them to be widely adopted [133][135].

New challenges arise from applying neural compression, or deep learning in general, to new types of data. Often, finding an effective digital representation of the data that meets application requirements is an art in itself, such as complex 3D scenes that need to be rendered from arbitrary viewing angles [129][212][136]. The next, and often highly related question, is how to best process and ultimately compress such data with neural networks, e.g., in the non-linear transform coding paradigm. For example, point clouds can be compressed using voxelization followed by non-linear transform coding [147][72], but a naive implementation using a 3D convolutional autoencoder can quickly run out of memory, and care is needed to apply neural networks strategically [72][148]. In the case of image data, these basic issues around representation and neural architectures have been well addressed prior to the current wave of research in neural image compression: modern computers have long been able to efficiently represent and manipulate digital images as matrices, and the machine learning community have converged to convolution neural networks as the default architecture for image processing. For many emerging data types, these issues have yet to be fully resolved, and the solutions may involve substantial domain knowledge and well-designed data structures [212][136]. We refer interested readers to the survey by Quach *et al.* [146] for a more in-depth discussion of these challenges in the context of point cloud compression.

Many open questions also remain in the design of objective functions in lossy compression. For example, it is poorly understood to which extent divergences and adversarial approaches are needed for realism, or

whether realism could also be achieved through well-crafted distortions and no-reference metrics. Adversarial losses also continue to pose challenges for tuning and optimization and no loss has yet emerged which can be trusted to reliably judge the perceptual quality of reconstructions in training and evaluation. On a related note, it remains to be seen how well neural networks can optimize various compression objectives compared to theoretical performance limits [211][199], and to better understand the cause of suboptimality.

The use of quantization continues to cause a mismatch between training and test time performance, and how much this affects compression performance is still not clearly understood. Reverse channel coding is a promising alternative which eliminates the need for quantization, but has only recently been considered in neural compression [186]. Open questions include the design of efficient coding schemes and the impact these schemes have on performance when compared to approaches based on quantization.

With the advances in neural lossy compression also emerges the risk of miscommunication, especially in semantically constrained domains. The reconstructions from neural lossy compression models, especially those targeting realism [152][3][127][185], can appear highly realistic at extremely low bit-rates, yet misrepresent the semantic content or other “relevant” information in the original data. The reconstructions might also be stochastically generated, and differ across different users and times of access. These effects raise safety and ethics concerns where such miscommunication can have severe consequences, e.g., in the transmission of security camera videos. Besides building safeguards into our methods, addressing these concerns may also require re-examining our choice of objective functions, for example, the distortion function and its effectiveness at capturing what is truly “relevant” to the end user of lossy compression.

Acknowledgements

The authors thank Karen Ullrich, Yingzhen Li, Thanasi Bakis, and David Minnen for providing valuable feedback on our manuscript. Yibo Yang acknowledges support from the Hasso Plattner Research School at UC Irvine. Stephan Mandt acknowledges support by the National Science Foundation (NSF) under the NSF CAREER Award 2047418; NSF Grants 2003237 and 2007719, the Department of Energy, Office of Science under grant DE-SC0022331, as well as gifts from Intel, Disney, and Qualcomm.

References

- [1] E. Agustsson, F. Mentzer, M. Tschannen, L. Cavigelli, R. Timofte, L. Benini, and L. V. Gool, “Soft-to-Hard Vector Quantization for End-to-End Learning Compressible Representations,” in *Advances in Neural Information Processing Systems 30*, pp. 1141–1151, 2017.
- [2] E. Agustsson and L. Theis, “Universally Quantized Neural Compression,” in *Advances in Neural Information Processing Systems 33*, 2020.
- [3] E. Agustsson, M. Tschannen, F. Mentzer, R. Timofte, and L. V. Gool, “Generative adversarial networks for extreme learned image compression,” in *The IEEE International Conference on Computer Vision (ICCV)*, pp. 221–231, 2019.
- [4] E. Agustsson, D. Minnen, N. Johnston, J. Ballé, S. J. Hwang, and G. Toderici, “Scale-space flow for end-to-end optimized video compression,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8503–8512, 2020.
- [5] N. Ahmed, T. R. Natarajan, and K. R. Rao, “Discrete cosine transform,” *IEEE Transactions on Computers*, vol. C-23, 1974, pp. 90–93.
- [6] A. Alemi, B. Poole, I. Fischer, J. Dillon, R. A. Saurous, and K. Murphy, “Fixing a broken elbo,” in *International Conference on Machine Learning*, PMLR, pp. 159–168, 2018.

- [7] A. A. Alemi, I. Fischer, J. V. Dillon, and K. Murphy, “Deep variational information bottleneck,” *arXiv preprint arXiv:1612.00410*, 2016.
- [8] D. Amir and Y. Weiss, “Understanding and simplifying perceptual distances,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12 226–12 235, Jun. 2021.
- [9] S. Arimoto, “An algorithm for computing the capacity of arbitrary discrete memoryless channels,” *IEEE Transactions on Information Theory*, vol. 18, no. 1, 1972, pp. 14–20.
- [10] J. Ballé, P. A. Chou, D. Minnen, S. Singh, N. Johnston, E. Agustsson, S. J. Hwang, and G. Toderici, “Nonlinear transform coding,” *IEEE Trans. on Special Topics in Signal Processing*, vol. 15, 2021.
- [11] J. Ballé, V. Laparra, and E. P. Simoncelli, “End-to-end optimization of nonlinear transform codes for perceptual quality,” in *Picture Coding Symposium*, 2016.
- [12] J. Ballé, V. Laparra, and E. P. Simoncelli, “End-to-end Optimized Image Compression,” in *International Conference on Learning Representations*, 2017.
- [13] J. Ballé, N. Johnston, and D. Minnen, “Integer networks for data compression with latent-variable models,” in *International Conference on Learning Representations*, 2018.
- [14] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston, “Variational Image Compression with a Scale Hyperprior,” *International Conference on Learning Representations*, 2018.
- [15] R. Bamler, “Understanding Entropy Coding With Asymmetric Numeral Systems (ANS): a Statistician’s Perspective,” *arXiv preprint arXiv:2201.01741*, 2022.
- [16] F. Bellard, *Lossless Data Compression with Neural Networks*, 2019.
- [17] Y. Bengio, N. Leonard, and A. Courville, *Estimating or propagating gradients through stochastic neurons for conditional computation*, 2013.

- [18] C. H. Bennett and P. W. Shor, “Entanglement-Assisted Capacity of a Quantum Channel and the Reverse Shannon Theorem,” *IEEE Transactions on Information Theory*, vol. 48, no. 10, 2002.
- [19] T. Berger and J. D. Gibson, “Lossy source coding,” *IEEE Transactions on Information Theory*, vol. 44, no. 6, 1998, pp. 2693–2723.
- [20] T. G. Bever and D. Poeppel, “Analysis by synthesis: A (re-) emerging program of research for language and vision,” *Biolinguistics*, vol. 4, no. 2-3, 2010, pp. 174–200.
- [21] S. Bhardwaj, I. Fischer, J. Ballé, and T. Chinen, “An Unsupervised Information-Theoretic Perceptual Quality Metric,” in *Advances in Neural Information Processing Systems*, vol. 33, pp. 13–24, Curran Associates, Inc., 2020.
- [22] T. Bird, J. Ballé, S. Singh, and P. A. Chou, “3d scene compression through entropy penalized neural representation functions,” in *2021 Picture Coding Symposium (PCS)*, IEEE, pp. 1–5, 2021.
- [23] C. M. Bishop, *Pattern Recognition and Machine Learning*, M. Jordan, J. Kleinberg, and B. Schölkopf, Eds., ser. Information science and statistics. Springer, 2006, ch. Graphical. DOI: [10.1117/1.2819119](https://doi.org/10.1117/1.2819119).
- [24] R. Blahut, “Computation of channel capacity and rate-distortion functions,” *IEEE Transactions on Information Theory*, vol. 18, no. 4, 1972, pp. 460–473. DOI: [10.1109/TIT.1972.1054855](https://doi.org/10.1109/TIT.1972.1054855).
- [25] Y. Blau and T. Michaeli, “The perception-distortion tradeoff,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6228–6237, 2018.
- [26] G. E. Blelloch, *Introduction to Data Compression*. Carnegie Mellon University, 2013.
- [27] L. Boltzmann, “Vorlesungen über Gastheorie,(2 volumes),” *Leipzig (1895, 1898)*, 1895.
- [28] S. Bosse, D. Maniry, T. Wiegand, and W. Samek, “A deep neural network for image quality assessment,” in *2016 IEEE International Conference on Image Processing (ICIP)*, pp. 3773–3777, 2016. DOI: [10.1109/ICIP.2016.7533065](https://doi.org/10.1109/ICIP.2016.7533065).

- [29] H. Bourlard and Y. Kamp, “Auto-association by multilayer perceptrons and singular value decomposition,” *Biological cybernetics*, vol. 59, no. 4, 1988, pp. 291–294.
- [30] J. Bruna, P. Sprechmann, and Y. LeCun, “Super-resolution with deep convolutional sufficient statistics,” in *International Conference on Learning Representations*, Jan. 2016.
- [31] J. Bruna and S. Mallat, “Invariant scattering convolution networks,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, 2013, pp. 1872–1886. DOI: [10.1109/TPAMI.2012.230](https://doi.org/10.1109/TPAMI.2012.230).
- [32] Y. Burda, R. Grosse, and R. Salakhutdinov, “Importance weighted autoencoders,” *arXiv preprint arXiv:1509.00519*, 2015.
- [33] T. Chen, H. Liu, Q. Shen, T. Yue, X. Cao, and Z. Ma, “Deepcoder: A deep neural network based video compression,” in *2017 IEEE Visual Communications and Image Processing (VCIP)*, IEEE, pp. 1–4, 2017.
- [34] X. Chen, N. Mishra, M. Rohaninejad, and P. Abbeel, “PixelSNAIL: An improved autoregressive generative model,” in *Proceedings of the 35th International Conference on Machine Learning*, J. Dy and A. Krause, Eds., pp. 864–872, PMLR, Jul. 2018.
- [35] Z. Cheng, H. Sun, M. Takeuchi, and J. Katto, “Learned image compression with discretized gaussian mixture likelihoods and attention modules,” *arXiv preprint arXiv:2001.01568*, 2020. arXiv: [2001.01568](https://arxiv.org/abs/2001.01568) [eess.IV].
- [36] Y. Choi, M. El-Khamy, and J. Lee, “Variable rate deep image compression with a conditional autoencoder,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2019.
- [37] P. A. Chou, T. Lookabaugh, and R. M. Gray, “Entropy-constrained vector quantization,” *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 37, no. 1, 1989, pp. 31–42.
- [38] R. Clausius, *Mechanical Theory of Heat*. Taylor and Francis, 1850.
- [39] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, vol. 2. John Wiley & Sons, 2006.

- [40] C. Cremer, X. Li, and D. Duvenaud, “Inference suboptimality in variational autoencoders,” in *International Conference on Machine Learning*, pp. 1078–1086, 2018.
- [41] P. Cuff, “Communication requirements for generating correlated random variables,” in *2008 IEEE International Symposium on Information Theory*, pp. 1393–1397, 2008.
- [42] P. W. Cuff and E. C. Song, “The likelihood encoder for source coding,” in *2013 IEEE Information Theory Workshop*, 2013.
- [43] J. Degraeve and I. Korshunova, *How we can make machine learning algorithms tunable*. URL: <https://www.engraved.blog/how-we-can-make-machine-learning-algorithms-tunable/>.
- [44] E. L. Denton, S. Chintala, A. Szlam, and R. Fergus, “Deep generative image models using a Laplacian pyramid of adversarial networks,” in *Advances in Neural Information Processing Systems (NIPS)*, pp. 1486–1494, 2015.
- [45] K. Ding, K. Ma, S. Wang, and E. P. Simoncelli, “Comparison of Full-Reference Image Quality Models for Optimization of Image Processing Systems,” *International Journal of Computer Vision*, no. 129, 2021, pp. 1258–1281.
- [46] L. Dinh, D. Krueger, and Y. Bengio, “Nice: Non-linear independent components estimation,” *arXiv preprint arXiv:1410.8516*, 2014.
- [47] L. Dinh, J. Sohl-Dickstein, and S. Bengio, “Density estimation using real nvp,” *arXiv preprint arXiv:1605.08803*, 2016.
- [48] A. Djelouah, J. Campos, S. Schaub-Meyer, and C. Schroers, “Neural inter-frame compression for video coding,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6421–6429, 2019.
- [49] C. Dong, C. C. Loy, K. He, and X. Tang, “Image super-resolution using deep convolutional networks,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 2, 2015, pp. 295–307.

- [50] A. Dosovitskiy and T. Brox, “Generating images with perceptual similarity metrics based on deep networks,” in *Advances in Neural Information Processing Systems*, vol. 29, Curran Associates, Inc., 2016. URL: <https://proceedings.neurips.cc/paper/2016/file/371bce7dc83817b7893bcdeed13799b5-Paper.pdf>.
- [51] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. Van Der Smagt, D. Cremers, and T. Brox, “Flownet: Learning optical flow with convolutional networks,” in *Proceedings of the IEEE international conference on computer vision*, pp. 2758–2766, 2015.
- [52] Z. Duan, M. Lu, Z. Ma, and F. Zhu, “Lossy image compression with quantized hierarchical vaes,” *arXiv preprint arXiv:2208.F180313056*, 2022.
- [53] Y. Dubois, B. Bloem-Reddy, K. Ullrich, and C. J. Maddison, “Lossy compression for lossless prediction,” in *Neural Information Processing Systems*, 2021.
- [54] J. Duda, “Asymmetric numeral systems,” *arXiv preprint arXiv:0902.0271*, 2009.
- [55] A. M. Eskicioglu and P. S. Fisher, “Image quality measures and their performance,” *IEEE Trans. Communications*, vol. 43, 1995, pp. 2959–2965.
- [56] P. Esser, R. Rombach, and B. Ommer, “Taming transformers for high-resolution image synthesis,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12 873–12 883, 2021.
- [57] G. Flamich, M. Havasi, and J. M. Hernández-Lobato, *Compressing Images by Encoding Their Latent Representations with Relative Entropy Coding*, 2020.
- [58] G. Flamich, S. Markou, and J. M. Hernandez-Lobato, “Fast relative entropy coding with a* coding,” in *Proceedings of the 39th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 162, pp. 6548–6577, PMLR, Jul. 2022.
- [59] B. J. Frey, *Bayesian networks for pattern classification, data compression, and channel coding*. Citeseer, 1998.

- [60] B. J. Frey and G. E. Hinton, "Efficient stochastic source coding and an application to a bayesian network source model," *The Computer Journal*, vol. 40, no. 2_and_3, 1997, pp. 157–165.
- [61] S. M. G.L. Sicuranza G. Ramponi, "Artificial neural network for image compression," *Electronics Letters*, vol. 26, 7 Mar. 1990, 477–479(2).
- [62] L. A. Gatys, A. S. Ecker, and M. Bethge, "Image style transfer using convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Jun. 2016. URL: http://www.cv-foundation.org/openaccess/content_cvpr_2016/html/Gatys_Image_Style_Transfer_CVPR_2016_paper.html.
- [63] A. Gersho and R. M. Gray, *Vector quantization and signal compression*, vol. 159. Springer Science & Business Media, 2012.
- [64] A. Golinski, R. Pourreza, Y. Yang, G. Sautiere, and T. S. Cohen, "Feedback recurrent autoencoder for video compression," in *Proceedings of the Asian Conference on Computer Vision*, 2020.
- [65] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative Adversarial Nets," in *Advances in Neural Information Processing Systems*, Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger, Eds., vol. 27, 2014.
- [66] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- [67] V. K. Goyal, "Theoretical foundations of transform coding," *IEEE Signal Processing Magazine*, vol. 18, no. 5, 2001, pp. 9–21.
- [68] V. K. Goyal, J. Zhuang, and M. Veiterli, "Transform coding with backward adaptive updates," *IEEE Transactions on Information Theory*, vol. 46, no. 4, 2000, pp. 1623–1633.
- [69] R. M. Gray and D. L. Neuhoff, "Quantization," *IEEE transactions on information theory*, vol. 44, no. 6, 1998, pp. 2325–2383.
- [70] K. Gregor, F. Besse, D. Jimenez Rezende, I. Danihelka, and D. Wierstra, "Towards conceptual compression," *Advances In Neural Information Processing Systems*, vol. 29, 2016, pp. 3549–3557.

- [71] P. D. Grünwald, *The minimum description length principle*. MIT press, 2007.
- [72] A. F. Guarda, N. M. Rodrigues, and F. Pereira, “Point cloud coding: Adopting a deep learning-based approach,” in *2019 Picture Coding Symposium (PCS)*, IEEE, pp. 1–5, 2019.
- [73] A. Habibian, T. v. Rozendaal, J. M. Tomczak, and T. S. Cohen, “Video compression with rate-distortion autoencoders,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7033–7042, 2019.
- [74] J. Han, S. Lombardo, C. Schroers, and S. Mandt, “Deep generative video compression,” in *Neural Information Processing Systems*, 2019.
- [75] P. Harsha, R. Jain, D. McAllester, and J. Radhakrishnan, “The Communication Complexity of Correlation,” in *Twenty-Second Annual IEEE Conference on Computational Complexity*, pp. 10–23, 2007.
- [76] M. Havasi, R. Peharz, and J. M. Hernández-Lobato, “Minimal Random Code Learning: Getting Bits Back from Compressed Model Parameters,” in *International Conference on Learning Representations*, 2019.
- [77] L. Helminger, A. Djelouah, M. Gross, and C. Schroers, “Lossy image compression with normalizing flows,” *arXiv preprint arXiv:2008.10486*, 2020.
- [78] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, “GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium,” in *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [79] A. Hines and N. Harte, “Speech Intelligibility Prediction using a Neurogram Similarity Index Measure,” *Speech Communication*, vol. 54, no. 2, 2012, pp. 306–320.
- [80] A. Hines, J. Skoglund, A. C. Kokaram, and N. Harte, “ViSQOL: an objective speech quality model,” *EURASIP Journal on Audio, Speech, and Music Processing*, 2015. DOI: [10.1186/s13636-015-0054-9](https://doi.org/10.1186/s13636-015-0054-9).

- [81] G. E. Hinton and D. Van Camp, “Keeping the neural networks simple by minimizing the description length of the weights,” in *Proceedings of the sixth annual conference on Computational learning theory*, pp. 5–13, 1993.
- [82] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” in *NeurIPS*, 2020.
- [83] J. Ho, E. Lohn, and P. Abbeel, “Compression with flows via local bits-back coding,” in *Advances in Neural Information Processing Systems*, pp. 3874–3883, 2019.
- [84] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, vol. 9, no. 8, 1997, pp. 1745–1780.
- [85] A. Honkela and H. Valpola, “Variational learning and bits-back coding: An information-theoretic view to bayesian learning,” *IEEE transactions on Neural Networks*, vol. 15, no. 4, 2004, pp. 800–810.
- [86] E. Hoogeboom, A. A. Gritsenko, J. Bastings, B. Poole, R. v. d. Berg, and T. Salimans, “Autoregressive diffusion models,” *arXiv preprint arXiv:2110.02037*, 2021.
- [87] E. Hoogeboom, J. Peters, R. van den Berg, and M. Welling, “Integer discrete flows and lossless compression,” in *Advances in Neural Information Processing Systems*, pp. 12 134–12 144, 2019.
- [88] R. Hosseini, F. Sinz, and M. Bethge, “Lower bounds on the redundancy of natural images,” *Vision Research*, vol. 50, no. 22, 2010, pp. 2213–2222.
- [89] L. Huang, S. Wang, K. Wong, J. Liu, and R. Urtasun, “Oct-squeeze: Octree-structured entropy model for lidar compression,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 1313–1323, 2020.
- [90] G. Hudson, A. Léger, B. Niss, I. Sebestyén, and J. Vaaben, “JPEG-1 standard 25 years: Past, present, and future reasons for a success,” *Journal of Electronic Imaging*, vol. 27, no. 4, 2018, p. 040 901.
- [91] D. A. Huffman, “A method for the construction of minimum-redundancy codes,” *Proceedings of the IRE*, vol. 40, no. 9, 1952, pp. 1098–1101.

- [92] ITU-T, *Recommendation ITU-T T.81: Information technology – Digital compression and coding of continuous-tone still images – Requirements and guidelines*, 1992.
- [93] ITU-T, *Recommendation ITU-T T.800.2: Methods for objective and subjective assessment of speech and video quality*, 2016.
- [94] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton, “Adaptive mixtures of local experts,” *Neural Computation*, 1991.
- [95] E. Jang, S. Gu, and B. Poole, “Categorical reparameterization with gumbel-softmax,” *arXiv preprint arXiv:1611.01144*, 2016.
- [96] J. Jiang, “Image compression with neural networks—a survey,” *Signal processing: image Communication*, vol. 14, no. 9, 1999, pp. 737–760.
- [97] J. Kim, A.-D. Nguyen, and S. Lee, “Deep cnn-based blind image quality predictor,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 1, 2019, pp. 11–24. DOI: [10.1109/TNNLS.2018.2829819](https://doi.org/10.1109/TNNLS.2018.2829819).
- [98] Y. Kim, S. Wiseman, A. Miller, D. Sontag, and A. Rush, “Semi-amortized variational autoencoders,” in *International Conference on Machine Learning*, PMLR, pp. 2678–2687, 2018.
- [99] D. Kingma and M. Welling, “Auto-encoding variational Bayes,” in *International Conference on Learning Representations*, 2014.
- [100] D. P. Kingma and P. Dhariwal, “Glow: Generative Flow with Invertible 1x1 Convolutions,” in *Advances in Neural Information Processing Systems 31*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., pp. 10 215–10 224, 2018.
- [101] D. P. Kingma, T. Salimans, R. Jozefowicz, X. Chen, I. Sutskever, and M. Welling, “Improved variational inference with inverse autoregressive flow,” *Advances in neural information processing systems*, vol. 29, 2016.
- [102] B. Knoll, *Cmix*. URL: <https://www.byronknoll.com/cmixon.html>.

- [103] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds., vol. 25, Curran Associates, Inc., 2012. URL: <https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>.
- [104] H. Larochelle and I. Murray, “The Neural Autoregressive Distribution Estimator,” *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics*, 2011.
- [105] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi, “Photo-realistic single image super-resolution using a generative adversarial network,” in *CVPR*, 2017.
- [106] J. Lee, S. Cho, and S.-K. Beack, “Context-adaptive entropy model for end-to-end optimized image compression,” in *International Conference on Learning Representations*, May 2019.
- [107] C. T. Li and A. E. Gamal, “Strong Functional Representation Lemma and Applications to Coding Theorems,” in *IEEE International Symposium on Information Theory*, 2017.
- [108] M. Li, W. Zuo, S. Gu, D. Zhao, and D. Zhang, *Learning convolutional networks for content-weighted image compression*, 2017. arXiv: [1703.10553](https://arxiv.org/abs/1703.10553) [cs.CV].
- [109] Z. Li, A. Aaron, I. Katsavounidis, A. Moorthy, and M. Manohara, *Toward a practical perceptual video quality metric*, 2016. URL: <https://netflixtechblog.com/toward-a-practical-perceptual-video-quality-metric-653f208b9652>.
- [110] A. Liu, S. Mandt, and G. V. d. Broeck, “Lossless compression with probabilistic circuits,” in *International Conference on Learning Representations*, 2022.
- [111] D. Liu, Y. Li, J. Lin, H. Li, and F. Wu, “Deep learning-based video coding: A review and a case study,” *ACM Computing Surveys (CSUR)*, vol. 53, no. 1, 2020, pp. 1–35.
- [112] H. Liu, T. Chen, P. Guo, Q. Shen, X. Cao, Y. Wang, and Z. Ma, “Non-local attention optimized deep image compression,” *arXiv preprint arXiv:1904.09757*, 2019.

- [113] T.-J. Liu, Y.-C. Lin, W. Lin, and C.-C. J. Kuo, “Visual quality assessment: Recent developments, coding applications and future trends,” *APSIPA Transactions on Signal and Information Processing*, vol. 2, 2013. DOI: [10.1017/ATSIP.2013.5](https://doi.org/10.1017/ATSIP.2013.5).
- [114] S. Lloyd, “Least squares quantization in pcm,” *IEEE transactions on information theory*, vol. 28, no. 2, 1982, pp. 129–137.
- [115] G. Lu, W. Ouyang, D. Xu, X. Zhang, C. Cai, and Z. Gao, “Dvc: An end-to-end deep video compression framework,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11 006–11 015, 2019.
- [116] G. Lu, X. Zhang, W. Ouyang, L. Chen, Z. Gao, and D. Xu, “An end-to-end learning framework for video compression,” *IEEE transactions on pattern analysis and machine intelligence*, 2020.
- [117] S. Ma, X. Zhang, C. Jia, Z. Zhao, S. Wang, and S. Wang, “Image and video compression with neural networks: A review,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 6, 2019, pp. 1683–1698.
- [118] D. J. MacKay and D. J. Mac Kay, *Information theory, inference and learning algorithms*. Cambridge university press, 2003.
- [119] C. J. Maddison, A. Mnih, and Y. W. Teh, “The concrete distribution: A continuous relaxation of discrete random variables,” *arXiv preprint arXiv:1611.00712*, 2016.
- [120] M. Mahoney, *Large Text Compression Benchmark*. URL: <http://mattmahoney.net/dc/text.html>.
- [121] X. Mao, Q. Li, H. Xie, R. Y. K. Lau, Z. Wang, and S. P. Smolley, “Least squares generative adversarial networks,” in *IEEE International Conference on Computer Vision (ICCV)*, pp. 2813–2821, 2017. DOI: [10.1109/ICCV.2017.304](https://doi.org/10.1109/ICCV.2017.304).
- [122] G. Martin, “Range encoding: An algorithm for removing redundancy from a digitised message,” in *Video and Data Recording Conference, Southampton, 1979*, pp. 24–27, 1979.
- [123] Y. Matsubara, R. Yang, M. Levorato, and S. Mandt, “Supervised compression for resource-constrained edge computing systems,” *arXiv preprint arXiv:2108.11898*, 2021.

- [124] F. Mentzer, E. Agustsson, M. Tschannen, R. Timofte, and L. Van Gool, “Conditional probability models for deep image compression,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4394–4402, 2018.
- [125] F. Mentzer, E. Agustsson, M. Tschannen, R. Timofte, and L. Van Gool, “Practical full resolution learned lossless image compression,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [126] F. Mentzer, G. Toderici, D. Minnen, S.-J. Hwang, S. Caelles, M. Lucic, and E. Agustsson, “Vct: A video compression transformer,” *arXiv preprint arXiv:2206.07307*, 2022.
- [127] F. Mentzer, G. D. Toderici, M. Tschannen, and E. Agustsson, “High-fidelity generative image compression,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [128] B. Meyer and P. Tischer, “Glicbawls – Grey Level Image Compression By Adaptive Weighted Least Squares,” in *Data Compression Conference*, 2001.
- [129] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, “Nerf: Representing scenes as neural radiance fields for view synthesis,” *Communications of the ACM*, vol. 65, no. 1, 2021, pp. 99–106.
- [130] T. Minka, “Divergence measures and message passing,” Microsoft Research, Tech. Rep. TR-2005-173, 2005.
- [131] D. Minnen, J. Ballé, and G. D. Toderici, “Joint Autoregressive and Hierarchical Priors for Learned Image Compression,” in *Advances in Neural Information Processing Systems 31*, 2018.
- [132] D. Minnen and S. Singh, “Channel-wise autoregressive entropy models for learned image compression,” in *IEEE International Conference on Image Processing (ICIP)*, 2020.
- [133] D. Minnen, *Current Frontiers In Neural Image Compression: The Rate-Distortion-Computation Trade-Off And Optimizing For Subjective Visual Quality*, Plenary talk at ICIP 2021, 2021. URL: <https://www.youtube.com/watch?v=84XkOIBhOas> (accessed on 10/26/2022).

- [134] A. Mittal, R. Soundararajan, and A. C. Bovik, “Making a “Completely Blind” Image Quality Analyzer,” *IEEE Signal Processing Letters*, vol. 20, no. 3, Mar. 2013, pp. 209–212. DOI: [10.1109/LSP.2012.2227726](https://doi.org/10.1109/LSP.2012.2227726).
- [135] D. Mukherjee, “Challenges in incorporating ML in a mainstream nextgen video codec,” *CLIC Workshop and Challenge on Learned Image Compression*, 2022. URL: https://storage.googleapis.com/clic2022_public/slides/Challenges%20in%20incorporating%20ML%20in%20a%20practical%20Nextgen%20Video%20Codec.pdf.
- [136] T. Müller, A. Evans, C. Schied, and A. Keller, “Instant neural graphics primitives with a multiresolution hash encoding,” *arXiv preprint arXiv:2201.05989*, 2022.
- [137] X. L. Nguyen, M. J. Wainwright, and M. I. Jordan, “On surrogate loss functions and f-divergences,” *The Annals of Statistics*, vol. 37, no. 2, 2009, pp. 876–904. DOI: [10.1214/08-AOS595](https://doi.org/10.1214/08-AOS595).
- [138] S. Nowozin, B. Cseke, and R. Tomioka, “f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization,” in *Advances in Neural Information Processing Systems*, vol. 29, 2016.
- [139] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. W. Senior, and K. Kavukcuoglu, “WaveNet: A Generative Model for Raw Audio,” *CoRR*, vol. abs/1609.03499, 2016. URL: <http://arxiv.org/abs/1609.03499>.
- [140] A. van den Oord and N. Kalchbrenner, “Pixel Recurrent Neural Networks,” in *Proceedings of the 33rd International Conference on Machine Learning*, 2016.
- [141] A. van den Oord, N. Kalchbrenner, O. Vinyals, A. G. L. Espeholt, and K. Kavukcuoglu, “Conditional Image Generation with PixelCNN Decoders,” in *Advances in Neural Information Processing Systems 29*, pp. 4790–4798, 2016.
- [142] A. v. d. Oord, O. Vinyals, and K. Kavukcuoglu, “Neural discrete representation learning,” 2017.

- [143] N. Parmar, A. Vaswani, J. Uszkoreit, L. Kaiser, N. Shazeer, A. Ku, and D. Tran, “Image transformer,” in *Proceedings of the 35th International Conference on Machine Learning*, J. Dy and A. Krause, Eds., ser. Proceedings of Machine Learning Research, vol. 80, pp. 4055–4064, Stockholmsmässan, Stockholm Sweden: PMLR, Jul. 2018.
- [144] R. Peharz, S. Lang, A. Vergari, K. Stelzner, A. Molina, M. Trapp, G. Van den Broeck, K. Kersting, and Z. Ghahramani, “Einsum networks: Fast and scalable learning of tractable probabilistic circuits,” in *International Conference on Machine Learning*, 2020.
- [145] J. C. Platt and A. H. Barr, “Constrained differential optimization for neural networks,” California Institute of Technology, Tech. Rep., 1988.
- [146] M. Quach, J. Pang, D. Tian, G. Valenzise, and F. Dufaux, “Survey on deep learning-based point cloud compression,” *Frontiers in Signal Processing*, 2022.
- [147] M. Quach, G. Valenzise, and F. Dufaux, “Learning convolutional transforms for lossy point cloud geometry compression,” in *2019 IEEE international conference on image processing (ICIP)*, IEEE, pp. 4320–4324, 2019.
- [148] M. Quach, G. Valenzise, and F. Dufaux, “Improved deep point cloud geometry compression,” in *2020 IEEE 22nd International Workshop on Multimedia Signal Processing (MMSP)*, IEEE, pp. 1–6, 2020.
- [149] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, “Zero-shot text-to-image generation,” in *International Conference on Machine Learning*, PMLR, pp. 8821–8831, 2021.
- [150] D. J. Rezende, S. Mohamed, and D. Wierstra, “Stochastic back-propagation and approximate inference in deep generative models,” in *International Conference on Machine Learning*, 2014.
- [151] D. Rezende and S. Mohamed, “Variational inference with normalizing flows,” in *International conference on machine learning*, PMLR, pp. 1530–1538, 2015.

- [152] O. Rippel and L. Bourdev, “Real-time adaptive image compression,” in *Proceedings of the 34th International Conference on Machine Learning*, 2017.
- [153] O. Rippel, S. Nair, C. Lew, S. Branson, A. G. Anderson, and L. Bourdev, “Learned video compression,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3454–3463, 2019.
- [154] J. Rissanen and G. G. Langdon, “Arithmetic coding,” *IBM Journal of research and development*, vol. 23, no. 2, 1979, pp. 149–162.
- [155] L. G. Roberts, “Picture Coding Using Pseudo-Random Noise,” *IRE Transactions on Information Theory*, 1962.
- [156] T. van Rozendaal, G. Sautière, and T. S. Cohen, *Lossy compression with distortion constrained optimization*, 2020. arXiv: [2005.04064](https://arxiv.org/abs/2005.04064) [cs.LG].
- [157] Y. Ruan, K. Ullrich, D. Severo, J. Townsend, A. Khisti, A. Doucet, A. Makhzani, and C. J. Maddison, “Improving lossless compression rates via monte carlo bits-back coding,” in *International Conference on Machine Learning*, 2021.
- [158] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Learning internal representations by error propagation,” California Univ San Diego La Jolla Inst for Cognitive Science, Tech. Rep., 1985.
- [159] C. Saharia, W. Chan, H. Chang, C. A. Lee, J. Ho, T. Salimans, D. J. Fleet, and M. Norouzi, “Palette: Image-to-Image Diffusion Models,” *CoRR*, vol. abs/2111.05826, 2021. arXiv: [2111.05826](https://arxiv.org/abs/2111.05826).
- [160] T. Salimans, A. Karpathy, X. Chen, and D. P. Kingma, “PixelCNN++: A pixelcnn implementation with discretized logistic mixture likelihood and other modifications,” in *International Conference on Learning Representations*, 2017.
- [161] S. Santurkar, D. Budden, and N. Shavit, “Generative compression,” in *2018 Picture Coding Symposium (PCS)*, pp. 258–262, 2018.

- [162] K. Sayood, “Vector quantization,” in *Introduction to Data Compression (Fourth Edition)*, ser. The Morgan Kaufmann Series in Multimedia Information and Systems, K. Sayood, Ed., 4th ed., Boston: Morgan Kaufmann, 2012, pp. 295–344. DOI: <https://doi.org/10.1016/B978-0-12-415796-5.00010-7>.
- [163] C. E. Shannon, “A Mathematical Theory of Communication,” *Bell System Technical Journal*, vol. 27, 1948, pp. 379–423.
- [164] C. Shannon, “Coding theorems for a discrete source with a fidelity criterion,” *IRE Nat. Conv. Rec., March 1959*, vol. 4, 1959, pp. 142–163.
- [165] H. Sheikh, M. Sabir, and A. Bovik, “A statistical evaluation of recent full reference image quality assessment algorithms,” *IEEE Transactions on Image Processing*, vol. 15, no. 11, 2006, pp. 3440–3451. DOI: [10.1109/TIP.2006.881959](https://doi.org/10.1109/TIP.2006.881959).
- [166] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *International Conference on Learning Representations*, 2015.
- [167] S. Singh, S. Abu-El-Haija, N. Johnston, J. Ballé, A. Shrivastava, and G. Toderici, “End-to-end learning of compressible features,” in *ICIP, IEEE*, pp. 3349–3353, 2020.
- [168] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, “Deep unsupervised learning using nonequilibrium thermodynamics,” in *International Conference on Machine Learning*, PMLR, pp. 2256–2265, 2015.
- [169] C. K. Sønderby, T. Raiko, L. Maaløe, S. K. Sønderby, and O. Winther, “Ladder variational autoencoders,” in *Advances in Neural Information Processing Systems*, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, Eds., vol. 29, Curran Associates, Inc., 2016. URL: <https://proceedings.neurips.cc/paper/2016/file/6ae07dcb33ec3b7c814df797cbda0f87-Paper.pdf>.
- [170] N. Sonehara, M. Kawato, S. Miyake, and K. Nakane, “Image data compression using a neural network model,” *International 1989 Joint Conference on Neural Networks*, 1989, 35–41 vol.2.

- [171] E. C. Song, P. Cuff, and H. V. Poor, “The likelihood encoder for lossy compression,” *IEEE Transactions on Information Theory*, vol. 62, 4 2016.
- [172] B. K. Sriperumbudur, K. Fukumizu, A. Gretton, B. Schölkopf, and G. R. G. Lanckriet, “On integral probability metrics, ϕ -divergences and binary classification,” *arXiv*, arXiv:0901.2698, Jan. 2009, arXiv:0901.2698.
- [173] M. Stern, N. Shazeer, and J. Uszkoreit, “Blockwise parallel decoding for deep autoregressive models,” in *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [174] K. Storrs, S. V. Leuven, S. Kojder, L. Theis, and F. Huszár, “Adaptive paired-comparison method for subjective video quality assessment on mobile devices,” in *Picture Coding Symposium*, 2018. URL: <https://arxiv.org/abs/1807.02175>.
- [175] M. Strathern, “Improving ratings: audit in the British University system,” *European Review*, vol. 5, 3 1997, pp. 305–321.
- [176] G. J. Sullivan, “Efficient scalar quantization of exponential and laplacian random variables,” *IEEE Transactions on Information Theory*, vol. 42, no. 5, 1996, pp. 1365–1374. DOI: [10.1109/18.532878](https://doi.org/10.1109/18.532878).
- [177] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” *arXiv*, vol. 1512.00567, 2015.
- [178] D. Tang, S. Singh, P. A. Chou, C. Hane, M. Dou, S. Fanello, J. Taylor, P. Davidson, O. G. Guleryuz, Y. Zhang, *et al.*, “Deep implicit volume compression,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 1293–1303, 2020.
- [179] L. Theis and E. Agustsson, “On the advantages of stochastic encoders,” 2021. URL: <https://arxiv.org/abs/2102.09270>.
- [180] L. Theis and M. Bethge, “Generative image modeling using spatial LSTMs,” in *Advances in Neural Information Processing Systems 28*, 2015.
- [181] L. Theis and J. Ho, “Importance weighted compression,” in *Neural Compression Workshop at ICLR 2021*, 2021. URL: https://openreview.net/forum?id=n6skss_9-v3.

- [182] L. Theis, R. Hosseini, and M. Bethge, “Mixtures of conditional Gaussian scale mixtures applied to multiscale image representations,” *PLoS ONE*, vol. 7, no. 7, 2012.
- [183] L. Theis, A. van den Oord, and M. Bethge, “A note on the evaluation of generative models,” in *International Conference on Learning Representations*, Apr. 2016. URL: <http://arxiv.org/abs/1511.01844>.
- [184] L. Theis, W. Shi, A. Cunningham, and F. Huszár, “Lossy Image Compression with Compressive Autoencoders,” in *International Conference on Learning Representations*, 2017.
- [185] L. Theis, T. Salimans, M. D. Hoffman, and F. Mentzer, *Lossy compression with gaussian diffusion*, 2022.
- [186] L. Theis and N. Yosri, “Algorithms for the communication of samples,” in *Proceedings of the 39th International Conference on Machine Learning*, 2021.
- [187] N. Tishby, F. C. Pereira, and W. Bialek, “The information bottleneck method,” *arXiv preprint physics/0004057*, 2000.
- [188] G. Toderici, S. M. O’Malley, S. J. Hwang, D. Vincent, D. Minnen, S. Baluja, M. Covell, and R. Sukthankar, “Variable rate image compression with recurrent neural networks,” in *International Conference on Learning Representations*, 2016.
- [189] G. Toderici, D. Vincent, N. Johnston, S. Jin Hwang, D. Minnen, J. Shor, and M. Covell, “Full resolution image compression with recurrent neural networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5306–5314, 2017.
- [190] R. Torfason, F. Mentzer, E. Agustsson, M. Tschannen, R. Timofte, and L. Van Gool, “Towards image understanding from deep compression without decoding,” *arXiv preprint arXiv:1803.06131*, 2018.
- [191] J. Townsend, “A tutorial on the range variant of asymmetric numeral systems,” *arXiv preprint arXiv:2001.09186*, 2020.
- [192] J. Townsend, T. Bird, J. Kunze, and D. Barber, “Hilloc: Lossless image compression with hierarchical latent variable models,” *arXiv preprint arXiv:1912.09953*, 2019.

- [193] J. Townsend, T. Bird, and D. Barber, “Practical lossless compression with latent variables using bits back coding,” *arXiv preprint arXiv:1901.04866*, 2019.
- [194] J. Townsend and I. Murray, “Lossless compression with state space models using bits back coding,” *arXiv preprint arXiv:2103.10150*, 2021.
- [195] D. Tran, K. Vafa, K. Agrawal, L. Dinh, and B. Poole, “Discrete flows: Invertible generative models of discrete data,” in *Advances in Neural Information Processing Systems*, pp. 14 719–14 728, 2019.
- [196] K. Tsubota and K. Aizawa, “Comprehensive comparisons of uniform quantizers for deep image compression,” in *2021 IEEE International Conference on Image Processing (ICIP)*, pp. 2089–2093, 2021. DOI: [10.1109/ICIP42928.2021.9506497](https://doi.org/10.1109/ICIP42928.2021.9506497).
- [197] B. Uria, I. Murray, and H. Larochelle, “RNADE: the real-valued neural autoregressive density-estimator,” in *Advances in Neural Information Processing Systems 26*, vol. 26, 2013.
- [198] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems*, vol. 30, pp. 5998–6008, 2017.
- [199] A. B. Wagner and J. Ballé, “Neural networks optimally compress the sawbridge,” in *2021 Data Compression Conference (DCC)*, IEEE, pp. 143–152, 2021.
- [200] C. S. Wallace, “Classification by minimum-message-length inference,” in *International Conference on Computing and Information*, Springer, pp. 72–81, 1990.
- [201] Z. Wang, E. P. Simoncelli, and A. C. Bovik, “Multiscale structural similarity for image quality assessment,” in *The Thrity-Seventh Asilomar Conference on Signals, Systems Computers, 2003*, vol. 2, 1398–1402 Vol.2, 2003. DOI: [10.1109/ACSSC.2003.1292216](https://doi.org/10.1109/ACSSC.2003.1292216).
- [202] R. J. Williams, “Simple statistical gradient-following algorithms for connectionist reinforcement learning,” *Machine learning*, vol. 8, no. 3, 1992, pp. 229–256.

- [203] C.-Y. Wu, N. Singhal, and P. Krahenbuhl, “Video compression through image interpolation,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 416–431, 2018.
- [204] X. Wu, K. U. Barthel, and W. Zhang, “Piecewise 2D Autoregression for Predictive Image Coding,” in *ICIP*, pp. 901–904, IEEE Computer Society, 1998.
- [205] Y. Xie, K. L. Cheng, and Q. Chen, “Enhanced invertible encoding for learned image compression,” in *Proceedings of the 29th ACM International Conference on Multimedia*, pp. 162–170, 2021.
- [206] R. Yang, F. Mentzer, L. Van Gool, and R. Timofte, “Learning for video compression with recurrent auto-encoder and recurrent probability model,” *IEEE Journal of Selected Topics in Signal Processing*, 2020.
- [207] R. Yang and S. Mandt, *Lossy image compression with conditional diffusion models*, 2022.
- [208] R. Yang, Y. Yang, J. Marino, and S. Mandt, “Hierarchical autoregressive modeling for neural video compression,” in *International Conference on Learning Representations*, 2020.
- [209] Y. Yang, R. Bamler, and S. Mandt, “Improving inference for neural image compression,” in *Neural Information Processing Systems (NeurIPS)*, 2020, 2020.
- [210] Y. Yang, R. Bamler, and S. Mandt, “Variational Bayesian Quantization,” in *International Conference on Machine Learning*, 2020.
- [211] Y. Yang and S. Mandt, “Towards empirical sandwich bounds on the rate-distortion function,” in *International Conference on Learning Representations*, 2022.
- [212] A. Yu, R. Li, M. Tancik, H. Li, R. Ng, and A. Kanazawa, “Plenotrees for real-time rendering of neural radiance fields,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5752–5761, 2021.
- [213] J. Yu, X. Li, J. Y. Koh, H. Zhang, R. Pang, J. Qin, A. Ku, Y. Xu, J. Baldridge, and Y. Wu, “Vector-quantized image modeling with improved vqgan,” *arXiv preprint arXiv:2110.04627*, 2021.

- [214] J. Yu, Y. Xu, J. Y. Koh, T. Luong, G. Baid, Z. Wang, V. Vasudevan, A. Ku, Y. Yang, B. K. Ayan, *et al.*, “Scaling autoregressive models for content-rich text-to-image generation,” *arXiv preprint arXiv:2206.10789*, 2022.
- [215] R. Zamir, *Lattice Coding for Signals and Networks*. Cambridge University Press, 2014.
- [216] R. Zamir and M. Feder, “On universal quantization by randomized uniform/lattice quantizers,” *IEEE Transactions on Information Theory*, vol. 38, no. 2, 1992, pp. 428–436.
- [217] C. Zhang, J. Bütepage, H. Kjellström, and S. Mandt, “Advances in variational inference,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 8, 2018, pp. 2008–2026.
- [218] L. Zhang, L. Zhang, and A. C. Bovik, “A feature-enriched completely blind image quality evaluator,” *IEEE Transactions on Image Processing*, vol. 24, no. 8, 2015, pp. 2579–2591. DOI: [10.1109/TIP.2015.2426416](https://doi.org/10.1109/TIP.2015.2426416).
- [219] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, “The Unreasonable Effectiveness of Deep Features as a Perceptual Metric,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 586–595, 2018. DOI: [10.1109/CVPR.2018.00068](https://doi.org/10.1109/CVPR.2018.00068).
- [220] L. Zhou, C. Cai, Y. Gao, S. Su, and J. Wu, “Variational Autoencoder for Low Bit-rate Image Compression,” in *Challenge on Learned Image Compression*, 2018.
- [221] Zhou Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, “Image quality assessment: From error visibility to structural similarity,” *IEEE Transactions on Image Processing*, vol. 13, no. 4, 2004, pp. 600–612. DOI: [10.1109/TIP.2003.819861](https://doi.org/10.1109/TIP.2003.819861).
- [222] Y. Zhu, Y. Yang, and T. Cohen, “Transformer-based transform coding,” in *International Conference on Learning Representations*, 2021.
- [223] J. Ziv, “On universal quantization,” *IEEE Transactions on Information Theory*, vol. 31, no. 3, 1985, pp. 344–347.
- [224] J. Ziv and A. Lempel, “A universal algorithm for sequential data compression,” *IEEE Transactions on information theory*, vol. 23, no. 3, 1977, pp. 337–343.