

LOSSLESS COMPRESSION WITH LATENT VARIABLE MODELS

James Townsend

A dissertation submitted in partial fulfillment
of the requirements for the degree of
Doctor of Philosophy
of
University College London.

Department of Computer Science
University College London

April 22, 2021

I, James Townsend, confirm that the work presented in this thesis is my own. Where information has been derived from other sources, I confirm that this has been indicated in the work.

Abstract

We develop a simple and elegant method for lossless compression using latent variable models, which we call ‘bits back with asymmetric numeral systems’ (BB-ANS). The method involves interleaving encode and decode steps, and achieves an optimal rate when compressing batches of data. We demonstrate it firstly on the MNIST test set, showing that state-of-the-art lossless compression is possible using a small variational autoencoder (VAE) model. We then make use of a novel empirical insight, that fully convolutional generative models, trained on small images, are able to generalize to images of arbitrary size, and extend BB-ANS to hierarchical latent variable models, enabling state-of-the-art lossless compression of full-size colour images from the ImageNet dataset. We describe ‘Craystack’, a modular software framework which we have developed for rapid prototyping of compression using deep generative models.

Impact Statement

The research on which this thesis is based has been published at two top tier machine learning conferences (Townsend et al., 2019; Townsend et al., 2020). It has also been included in the prestigious and highly popular ‘Deep Unsupervised Learning’ course, taught at UC Berkeley, taking up roughly one hour of lecture material. The lecture includes a number of verbatim extracts and diagrams from the paper Townsend et al. (2019), upon which Chapter 3 is based. At the time of writing, a video of the lecture had been viewed over 5 thousand times on YouTube¹. A number of recent publications by other researchers have been based upon our work, including Hogeboom et al. (2019), Ho et al. (2019), F. H. Kingma et al. (2019), and van den Berg et al. (2020).

There is great potential for translation of this work into improvements to widely used production compression systems, both in general purpose codecs which are used by billions of people on a daily basis as they browse the internet, and more bespoke systems for specific applications such as compression of medical or astronomical images. We have sought to maximise long-term impact by publishing the tutorial paper Townsend (2020), which is intended as an accessible introduction to the ideas in this thesis; as well as releasing open source software for easy prototyping of future ideas in the area which this thesis covers, namely lossless compression using deep generative models².

¹See youtu.be/0IoLKnAg6-s?t=5681.

²Software available at github.com/j-towns/craystack, mirrored at doi.org/10.5281/zenodo.4707276.

Acknowledgements

Thanks to Mark Herbster for encouraging me to start a PhD, without Mark I certainly would not have started on this journey, and I can't thank him enough for gently nudging me in this direction. Thanks to my co-authors Tom Bird and Julius Kunze, they are both gifted researchers and working with them has been an absolute joy. Thanks to my supervisor David Barber and to others who have given advice and feedback on this work, particularly Paul Rubenstein and Raza Habib. Thanks to the examiners, Matt Kusner and Iain Murray, for an exciting viva and excellent, thorough feedback. A million thanks to my family, to my partner Maud and to my best friend Anthony, for all of their love and support.

Preface

The connections between information theory and machine learning have long been known to be deep. The two fields are so closely related that they have been described as ‘two sides of the same coin’ by MacKay (2003), who insists, in the preface to his book, that they ‘belong together’. One particularly elegant connection is the essential equivalence between probabilistic models of data and lossless compression methods. The source coding theorem (Shannon, 1948) can be thought of as the fundamental theorem describing this idea, and Huffman coding (Huffman, 1952), arithmetic coding (AC; Witten et al., 1987) and the more recently developed asymmetric numeral systems (ANS; Duda, 2009) are actual algorithms for implementing lossless compression, given some kind of probabilistic model.

The field of machine learning has experienced an explosion of activity in recent years, and, amongst other things, we have seen major breakthroughs in probabilistic modelling of high dimensional data. Recurrent neural networks and autoregressive models based on masked convolution have been shown to be effective generative models for images, audio and natural language (Graves, 2014; van den Oord et al., 2016a; van den Oord et al., 2016b). These models are slow to sample from, at least in a naïve implementation, which means that decompression using these models is similarly slow³, however they do tend to achieve state of the art test log-likelihoods, and hence state of the art compression rates.

Another significant, recently invented type of probabilistic generative model is the variational autoencoder (VAE), first presented in D. P. Kingma and Welling (2014) and Rezende et al. (2014). VAEs are latent variable models which use a neural network for efficient posterior inference and where the generative model

³Sampling in autoregressive models based on masked convolution can be sped up drastically using dynamic programming (Le Paine et al., 2016; Ramachandran et al., 2017). However, there is not yet a general implementation of this method and at present a lot of developer time is required to hand-implement fast sampling for each individual model.

and posterior inference network are jointly trained, using stochastic gradient descent on the variational free energy, an objective function which we refer to in this work as the ‘evidence lower bound’ (ELBO). VAE models have been shown to obtain competitive (though usually not state of the art) log-likelihoods, and sampling from them tends to be much faster than it is from autoregressive and recurrent models.

In the last five years we have seen a number of papers covering applications of modern deep learning methods to *lossy* compression. Gregor et al. (2015) discusses applications of a VAE to compression, with an emphasis on lossy compression. Ballé et al. (2017), Theis et al. (2017), Ballé et al. (2018), and Minnen et al. (2018) all implement lossy compression using (variational) auto-encoder style models, and Tschannen et al. (2018) train a model for lossy compression using a GAN-like objective. Applications to *lossless* compression were less well covered in works published prior to Townsend et al. (2019), upon which Chapter 3 is based.

The classic lossless compression algorithms mentioned above (Huffman coding, AC, and ANS) do not naturally cater for latent variable models. However there is a method, known as ‘bits back coding’, which can be used to extend those algorithms to cope with such models (C. S. Wallace, 1990; Hinton and van Camp, 1993). Bits back coding was originally introduced as a purely theoretical construct, but was later implemented, in a primitive form, by Frey and Hinton (1996).

There is a fundamental incompatibility between the bits back method and the AC scheme upon which the primitive implementation was based. A workaround is suggested by Frey (1997), but this leads to a sub-optimal compression rate and an overly complex implementation. The central theoretical contribution of this thesis is a simple and elegant solution to the issue just mentioned, which involves implementing bits back using ANS instead of AC. We term the new coding scheme ‘Bits Back with ANS’ (BB-ANS).

Our scheme improves on early implementations of bits back coding in terms of compression rate and code complexity, allowing for efficient lossless compression of batches of data with deep latent variable models. We are also the first to implement a method, first suggested by MacKay (2003) for using bits back coding with continuous latent variables. In Chapter 3 we demonstrate the efficiency of BB-ANS by losslessly compressing the MNIST dataset with a VAE.

We find that BB-ANS with a VAE outperforms generic compression algorithms for both binarized and raw MNIST, with a very simple model architecture. In Chapter 4 we lay the ground work for scaling BB-ANS up to larger models, describing a method for vectorizing the underlying ANS coder, and discussing the generic software which we have written, enabling other machine learning researchers to easily prototype compression systems. In Chapter 5 we present an extension of BB-ANS to hierarchical latent variable models. We show that it is possible to use a fully convolutional model to compress images of arbitrary shape and size, and use this technique with BB-ANS to achieve a 6.5% improvement over the previous state of the art compression rate on full-size images from the ImageNet dataset.

In our experiments, we benchmark BB-ANS on image data because this is a well-studied domain for deep generative models. The (non-learned) image codecs which we benchmark against are Portable Network Graphics (PNG), WebP lossless, and Free Lossless Image Format⁴ (FLIF), which were first released in 1997, 2012, and 2015, respectively. All three of these codecs achieve compression using local, low-dimensional prediction. They adapt (effectively ‘learning’) as more pixels within a single image are processed. The broad approach taken in this thesis, which can be applied to data types other than images, and with models other than VAEs, is to spend a large amount of computation time adapting (training) a model on a generic dataset such as ImageNet, *before* it is presented with an image which it is tasked with compressing.

I would argue that the early results from using this approach, which are presented in this thesis, as well as the other recent works, are promising. In the period since the publication of Townsend et al. (2019), there have been a number of articles on this topic, most of which make use of ideas from our work. Mentzer et al. (2019) was published just after Townsend et al. (2019), and demonstrates learned lossless compression with a relatively weak model, achieving excellent run-times but failing to outperform existing methods in terms of compression rate. F. H. Kingma et al. (2019) builds directly on our work, proposing an extension of BB-ANS to hierarchical models; we compare this to our own extension in Chapter 5. Very recent work by Ruan et al. (2021) also directly extends the methods in this thesis, using Monte Carlo methods

⁴FLIF is being subsumed by the JPEG XL standard which is currently under development. JPEG XL is more general and achieves better lossless compression rates than FLIF, see Alakuijala et al., 2019.

to improve the compression rate. Hoogeboom et al. (2019), Ho et al. (2019), and van den Berg et al. (2020) all use ANS with flow models to do lossless compression. They show excellent lossless compression performance on 32×32 and 64×64 images but do not manage to scale their methods up to large images as we have. It is still an open question whether this potential paradigm shift can be exploited in the production-grade systems which are used by billions of people every day, and whether the gains will be as impressive in other domains, such as audio and video.

Structure of the thesis

This thesis is comprised of a background chapter (Chapter 1), followed by four ‘core’ chapters (Chapters 2 to 5). It is designed to be read by someone who has a little background knowledge in information theory and coding theory (MacKay, 2003, Chapters 4-6, provides an ideal introduction to the necessary topics), and who is already familiar with the types of modern generative models discussed above, particularly variational autoencoders (VAEs). A brief outline of the core chapters:

2. *Introduction to asymmetric numeral systems*

We introduce ANS coding, proving worst-case bounds on its performance. We include diagrams and pseudocode to assist the reader in their understanding. We also provide a 50 line working Python implementation to accompany this chapter, which can be found in Appendix A. This chapter is based on our tutorial paper, Townsend (2020). The ANS algorithm was first presented in Duda (2009).

3. *Bits back coding with asymmetric numeral systems*

We present our novel approach to bits-back coding and demonstrate its performance by compressing the MNIST test set with a small VAE model. We call this method ‘bits-back with asymmetric numeral systems’ (BB-ANS). We discuss a number of potential issues with the method and also discuss ways to improve this simple prototype system. This chapter is based on our paper, Townsend et al. (2019).

4. *Vectorizing ANS with Craystack*

We describe a method for implementing a vectorized ANS coder, drawing on earlier work by Giesen (2014). We describe ‘Craystack’, a software tool

which we have developed which aims to provide a flexible, user-friendly API for machine learning practitioners who wish to prototype compression systems. We discuss future directions for Craystack.

5. *Scaling up bits back coding with asymmetric numeral systems*

We demonstrate that the BB-ANS method can be scaled up to hierarchical VAEs and large, colour images from the ImageNet test set. The method achieves state of the art compression on a randomly selected subset of 2,000 images from the ImageNet dataset. We discuss the challenges and solutions which were necessary to achieve this scale-up. This chapter is based on our paper, Townsend et al. (2020).

Although each chapter in some ways builds on all of the previous, it should be possible to read and understand (at least on a high level) Chapter 4 without reading Chapter 3, and Chapter 5 without reading Chapter 4. Figure 1 visualizes this approximate dependency structure between the core chapters.

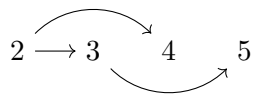


Figure 1: Dependencies between the core chapters of this thesis. We recommend not to read a chapter before reading its parents in this graph.

Contents

1	Background	18
1.1	Probabilistic generative models	18
1.1.1	Latent variable models	20
1.1.2	Further background on probabilistic generative models . .	22
1.2	Source coding	22
1.2.1	Terminology and the source coding theorem	22
1.2.2	Coding according to a model	23
1.3	Basic coding techniques	24
1.3.1	Huffman coding	24
1.3.2	Stream codes	25
1.4	Image compression in practice	28
1.4.1	Comparing PNG, FLIF, and our new codecs	29
2	Introduction to asymmetric numeral systems	32
2.1	Specifying the problem which ANS solves	33
2.2	Asymmetric numeral systems	34
2.2.1	The structure of the message	35
2.2.2	Constructing the pop operation	36
2.2.3	Popping in sequence	37
2.2.4	The function d	38
2.2.5	Computing s'	39
2.2.6	Pseudocode for d	41
2.2.7	Inverting the decoder	41
2.3	Further reading	43
3	Bits back coding with asymmetric numeral systems	44
3.1	Bits back coding	45
3.1.1	Bits back without ANS	45

3.1.2	Chaining bits back coding	47
3.1.3	Chaining bits back coding with ANS	47
3.2	Issues affecting the efficiency of BB-ANS	49
3.2.1	The extra information	50
3.2.2	Discretizing a continuous latent space	51
3.2.3	The need for ‘clean’ bits	52
3.3	Experiments	53
3.3.1	Using a VAE as the latent variable model	53
3.3.2	Compressing MNIST	53
3.4	Discussion	55
3.4.1	Extending BB-ANS to state-of-the-art latent variable models	55
3.4.2	Parallelization of BB-ANS	55
3.4.3	Communicating the model	56
3.5	Conclusion	56
4	Vectorizing ANS with Craystack	58
4.1	Vectorizing asymmetric numeral systems	58
4.1.1	The vectorized message data structure	58
4.1.2	Vectorizing the pop operation	59
4.1.3	The length of the vectorized message	60
4.1.4	The cost of vectorization	62
4.1.5	The benefit of vectorization	63
4.2	Craystack	63
4.2.1	High-level design	64
4.2.2	Generalizing vectorized ANS	64
4.3	Future directions for Craystack	65
5	Scaling up bits back coding with asymmetric numeral systems	67
5.1	Scaling up BB-ANS	68
5.1.1	Fully convolutional models	68
5.1.2	Reducing one-time overhead	68
5.1.3	Dynamic discretization	71
5.2	Experimental results	73
5.3	Discussion	75
5.4	Conclusion	76

<i>Contents</i>	13
6 General Conclusions	77
Appendices	78
A Example ANS implementation	78
B Reparameterizing discretized latents in hierarchical VAEs	81
C The ResNet VAE architecture	83
Bibliography	91

List of Figures

1	Dependencies between the core chapters of this thesis. We recommend not to read a chapter before reading its parents in this graph.	10
1.1	Comparing a photograph of the cat ‘Stasha’ with the compressed bit streams resulting from applying the codecs <code>gzip</code> , PNG and FLIF to the image. Vertical bands are clearly visible in the <code>gzip</code> stream, and upon close inspection there are some visible vertical bands in the PNG stream, implying some correlation between pixels and hence redundancy. No structure is visible in the FLIF stream, suggesting that FLIF may be close to optimal.	30
2.1	The two components of a message: the unsigned integer s (with $r_s = 16$) and the stack of unsigned integers t (with $r_t = 8$). The integers are represented here in base 2 (binary).	36
2.2	Showing the correspondence between s , $s \bmod 2^r$ and the symbol x . The interval $[0, 1] \subset \mathbb{R}$ is modelled by the set of integers $\{0, 1, \dots, 2^r - 1\}$. In this case $r = 3$ and the probabilities of each symbol are $P(\mathbf{a}) = 1/8$, $P(\mathbf{b}) = 2/8$, $P(\mathbf{c}) = 3/8$ and $P(\mathbf{d}) = 2/8$. . .	40
2.3	Showing the pmf of a distribution over symbols (left) and a visualization of the mapping d (middle and right). In the middle and right figures, numbers less than or equal to s_{\max} are plotted, for $s_{\max} = 20$ and $s_{\max} = 70$. The position of each number s plotted is set to the coordinates (x, s') , where $s', x = d(s)$. The heights of the bars are thus determined by the ratio s'/s from eq. (2.22), and can be seen to approach the heights of the lines in the histogram on the left (that is, to approach $P(x)$) as the density of numbers increases.	40

3.1	Visual comparison of 30 binarized MNIST images with bit stream outputs from running lossless compression algorithms PNG, bz2 and BB-ANS on the images.	45
3.2	Graphical model with latent variable z and observed variable x	46
3.3	Python implementation of BB-ANS encode ('push') and decode ('pop') methods.	50
3.4	A 2000 point moving average of the compression rate, in bits per dimension, during the compression process using BB-ANS with a VAE. We compress a concatenation of three shuffled copies of the MNIST test set.	55
4.1	The two components of a (vectorized) message: the vector of unsigned integers s (with $r_s = 16$) and the stack of unsigned integers t (with $r_t = 8$). The integers are represented here in base 2 (binary).	59
4.2	Runtime of vectorized vs. serial ANS implementations for different image sizes. Times were computed on a desktop with 6 CPU cores.	63
4.3	Visualizing the process of pushing images and latents from a VAE to the vectorized ANS stack with Craystack. The ANS stack head is shaped such that images and latents can be pushed and popped in parallel, without reshaping. Beneath the shaped top of the stack is the flat message stream output by ANS.	64
4.4	The BB-ANS combinator provided in Craystack. It accepts codecs for the prior, likelihood and approximate posterior as arguments and returns a codec which uses BB-ANS to compress data modelled with a latent variable model.	65
5.1	A selection of images from the ImageNet dataset and the compression rates achieved on the dataset by PNG, WebP, FLIF, Bit-Swap and the HiLLoC codec (with ResNet VAE) presented in this chapter.	68
5.2	Graphical models representing the generative and inference models with HiLLoC and Bit-Swap, both using a 3 layer latent hierarchy. The dashed lines indicate dependence on a fixed observation.	72

List of Tables

- 3.1 The shorthand we use for ANS encoding and decoding operations, based on left and right pointing arrows. The operations can be translated mechanically to or from the pseudocode shown in the right-hand column. 48
- 3.2 Visualizing the process by which a sender encodes (pushes) a symbol x onto an ANS message stack using Bits Back with ANS. The process starts with an existing ANS message, labelled ‘extra information’, and the operations in the right hand column are performed starting at the top of the table and working downwards. The ‘Variables’ column shows variables which are known to the sender before each operation is performed. 49
- 3.3 Compression rates on the binarized MNIST and full MNIST test sets, using BB-ANS and other benchmark compression schemes, measured in bits per dimension. Note that PNG and WebP are included for context but the comparison is not particularly fair, because the image files contain metadata, the size of which is significant relative to the size of an MNIST image, particularly in the binarized case. We also give (in brackets) the negative ELBO value for our trained VAEs, and the log-likelihood under the current state of the art model, ‘Locally masked PixelCNN’, all evaluated on the test set. 54
- 5.1 The Bit-Swap encoder and decoder procedures, in order from the top, for a model with L layers. For the encoder, the ‘Variables’ column shows variables known before each operation. For the decoder, that column shows variables known after each operation. 70

- 5.2 Compression performance of HiLLoC with RVAE compared to other codecs. Rates measured in bits/dimension (raw data is 8 bits/dimension). For HiLLoC we display compression rate and theoretical performance (ELBO). All HiLLoC results are obtained from the *same model*, trained on ImageNet 32. 74
- B.1 The BB-ANS encoding and decoding operations, in order from the top, for a hierarchical latent model with L layers. The Q_l are posterior distributions over the indices i_l of the discretized latent space for the l th latent, z_l . The discretization for the l th latent is created such that the intervals have equal mass under the prior. 82

Chapter 1

Background

This chapter aims to give a concise overview of the necessary background material for understanding the rest of the thesis and its context within the existing compression literature. We begin, in Section 1.1, by describing probabilistic generative models, and specifically latent variable models, which are a central tool in the compression methods developed in later chapters. Then in Section 1.2 we review some theoretical background material in compression, and in particular we define lossless compression and give a statement of the source coding theorem (Shannon, 1948). In Section 1.3, we outline a few of the most commonly used, generic, lossless compression methods. Finally, in Section 1.4, we motivate lossless compression of images and give an overview of two existing codecs, comparing their approach to ours.

1.1 Probabilistic generative models

A probabilistic model describes a variable (or set of variables) whose value is random. Such a model may be characterised by a probability ‘mass function’ $P: \mathcal{X} \rightarrow \mathbb{R}$, where \mathcal{X} is a countable (usually finite) set. The function P must satisfy two properties

1. $\forall x \in \mathcal{X}, \quad P(x) \geq 0$
2. $\sum_{x \in \mathcal{X}} P(x) = 1.$

We say that the distribution defined by P is ‘discrete’. This definition can be generalized to sets \mathcal{X} which are Lebesgue measurable (such as the set \mathbb{R} of real numbers). Then instead of a mass function the distribution is characterised by a ‘density function’ $p: \mathcal{X} \rightarrow \mathbb{R}$, with

1. $\forall x \in \mathcal{X}, \quad p(x) \geq 0$

$$2. \int_{x \in \mathcal{X}} p(x) = 1$$

and we say that the distribution is ‘continuous’, rather than discrete. As we will see in Chapter 2, we can efficiently compress outcomes from a discrete distribution if they can be broken down into a sequence of ‘symbols’, i.e. if we can write $x = x_1, x_2, \dots$ with $P(x) = \prod_t P(x_t | x_1, \dots, x_{t-1})$ and if we are able to compute the cumulative distribution function (CDF) and its inverse for each of the factors $P(x_t | x_1, \dots, x_{t-1})$.

In practice, a mass function or density function will usually have one or more tunable parameters. We usually use θ to denote the vector containing a model’s parameters and write $P(x; \theta)$ for the parametrized mass function. In the context of modern machine learning, the parameters of a probabilistic generative model are usually tuned (or ‘trained’), by optimizing the ‘log-likelihood’ function over a set of example data, referred to as the ‘training set’. That is, by solving the following optimization problem:

$$\hat{\theta} = \arg \max_{\theta} \sum_{x \in X} \log P(x; \theta), \quad (1.1)$$

where X is a set of example data. This technique is known as ‘maximum likelihood’ (ML) learning. If P can be tractably computed, then its gradient with respect to θ can usually be computed using automatic differentiation (AD), and the log-likelihood optimized using stochastic gradient descent (SGD). After optimization the model is usually evaluated on a held-out set of examples called the ‘test set’.

In the last five years, maximum likelihood training has been scaled to models with billions of parameters, and data sets with millions of examples, utilizing parallel graphics processing unit (GPU) based hardware for efficient training. Two particularly famous examples of systems which use this technique are WaveNet (van den Oord et al., 2016b), which is a probabilistic generative model for audio, used in Google’s speech synthesis applications, and GPT-3 (Brown et al., 2020), a natural language model with 175 billion parameters. Models for images which are trained using maximum likelihood are also fast approaching photo-realism in the samples which they generate. Recent examples include Menick and Kalchbrenner (2018), Jun et al. (2020), and Child (2020).

1.1.1 Latent variable models

Latent variable models are a class of probabilistic generative model which involve unobserved, or ‘latent’, variables. Their mass function is defined implicitly by an integral, or, to use the terminology of probability theory, by ‘marginalizing’ a latent variable:

$$P(x) := \int_z P(x|z)p(z)dz \quad (1.2)$$

The distribution specified by $p(z)$ is referred to as the ‘prior’, and the forward probability $P(x|z)$ the ‘likelihood’ (note likelihood here has a slightly different meaning to the definition in the previous section). We usually choose prior and likelihood distributions which are straightforward to sample from, and thus exact sampling from $P(x)$ is straightforward by first sampling z from the prior and then sampling x from the likelihood, conditioned on the sampled z .

A popular class of latent variable models, which we use to demonstrate our methods in Chapters 3 and 5, are called ‘variational auto-encoders’ (VAEs), first introduced in D. P. Kingma and Welling (2014) and Rezende et al. (2014). They usually use a simple, fixed prior distribution, such as a multivariate Gaussian with mean zero and with covariance equal to the identity matrix. The likelihood function $P(x|z)$ is usually of the form $P_{\text{simple}}(x|f(z;\theta))$, where f is a multi-layer neural network (i.e. a composition of differentiable, parametrized functions), and P_{simple} is a mass function which is straightforward to compute, and has the appropriate support. For discrete data, a discretized logistic distribution is often used (we give details of the distributions we used in our experiments in later chapters).

For VAE models the integral eq. (1.2) cannot easily be computed, and thus direct maximum likelihood training is not possible. To train a VAE, we instead optimize a variational lower bound on $\log P(x;\theta)$, called the ‘evidence lower-bound’ (ELBO), sometimes referred to as the ‘variational free energy’. It is defined as

$$\mathcal{L}(x;\theta,\phi) = \int_z q(z|x;\phi) \log \frac{P(x|z;\theta)p(z)}{q(z|x;\phi)} dz. \quad (1.3)$$

We refer to the newly introduced density function q as the variational posterior, or the approximate posterior, and the new parameters ϕ as the variational parameters. The fact that $\mathcal{L}(x;\theta,\phi) \leq \log P(x;\theta)$ follows directly from Jensen’s inequality.

The approximate posterior q can, in principal, be any distribution, but in VAEs it is common to use a parametrization similar to that of the likelihood, i.e. a function with the form $q_{\text{simple}}(z | g(x; \phi))$, where g is a multi-layer neural network and q_{simple} is usually a Gaussian distribution with a diagonal covariance matrix. Exact sampling from q is then straightforward, and this allows us to compute an unbiased Monte Carlo estimate of \mathcal{L} :

$$\hat{\mathcal{L}}(x; \theta, \phi) = \log \frac{P(x | z; \theta)p(z)}{q(z | x; \phi)} \text{ where } z \leftarrow q(\cdot | x; \phi). \quad (1.4)$$

The shorthand notation $z \leftarrow q(\cdot | x; \phi)$ means z is sampled from the approximate posterior distribution $q(z | x; \phi)$.

If the process used to generate z from q is differentiable with respect to ϕ , then the function $\hat{\mathcal{L}}$ can be differentiated with respect to θ and ϕ (we usually use automatic differentiation tools to do this), and SGD can be used to optimize \mathcal{L} .

Variational auto-encoders are reasonably straightforward to train and sample from, and since 2014 there have been a huge number of papers presenting different versions, with a general trend towards models with more parameters, better samples, and more accurate density estimation. Rather than survey this extensive literature we simply point to the two most recent examples of works which have pushed this envelope, Maaløe et al. (2019) and Vahdat and Kautz (2020), both of which demonstrate VAE image models with samples that, to the human eye, appear similar to real examples.

For the algorithms introduced later, we will need to be able to compress outcomes from the prior, likelihood and posterior of a latent variable model. As mentioned in Section 1.1, this means their mass functions must be factorizable into a product of conditionals in such a way that we are able to compute the CDF and its inverse under each conditional. Since, in a VAE, elements of the vectors x and z are usually modelled as independent under the three relevant distributions, we can simply use the natural factorization which has one factor for each element in the vector; it is also common practice to use distributions for which the CDF and inverse CDF for each element can easily be computed. Lossless compression is only possible for discrete (rather than continuous) data, which would seem to render it incompatible with the continuous latents typically used for VAEs. However, it turns out that it is straightforward to overcome this issue by quantizing continuous latents, in a way which has only a very

small effect on compression rates. We will give more detail in Chapter 3 and Chapter 5.

1.1.2 Further background on probabilistic generative models

For detailed background on probabilistic generative models, MacKay (2003) and Bishop (2006) are classic references. Murphy (2012) gives a more modern, extremely thorough, overview of the field, and the even more recent Goodfellow et al. (2016) includes VAE models and detail on the various neural network techniques which are useful for defining the functions $f(x; \theta)$ and $g(z; \phi)$ mentioned above.

1.2 Source coding

‘Source coding’, first formalized by Shannon (1948), is the name we give to the problem of trying to find lossless encodings for data which are drawn from a random source. Source coding is more-or-less synonymous with ‘lossless compression’. In this section we define basic terminology, then in Section 1.3 we review some of the basic algorithms for doing source coding. This is closely based on MacKay (2003) Chapters 4-7, which we recommend as an introduction to these topics.

1.2.1 Terminology and the source coding theorem

We use the following definition for probability distributions, based on that used by MacKay (2003):

Definition 1. An *ensemble* X is a triple $(x, \mathcal{A}_X, \mathcal{P}_X)$ where the *outcome* x is the value of a random variable, taking on one of a set of possible values $\mathcal{A}_X = \{a_1, \dots, a_I\}$, and $\mathcal{P}_X = \{p_1, \dots, p_I\}$ are non-negative real-valued probability weights with each $P(x = a_i) = p_i$ and therefore $\sum_{i=1}^I p_i = 1$.

The following are the basic quantities of concern in source coding and information theory in general. Here, and throughout the rest of the thesis, we use ‘log’ for the base 2 logarithm, usually denoted ‘log₂’.

Definition 2. The *Shannon information content* of an outcome x is

$$h(x = a_i) := \log \frac{1}{p_i}. \quad (1.5)$$

Definition 3. The *entropy* of an ensemble X is

$$H(X) := \sum_i p_i \log \frac{1}{p_i}. \quad (1.6)$$

The *source coding theorem* demonstrates the relevance of the entropy as a measure of a random variable's information content. We give an informal statement of the theorem. For more detail, including a proof of the theorem, see MacKay (2003), Chapter 4.

Theorem 1 (Source coding theorem, informal statement). *N i.i.d. random variables each with entropy $H(X)$ can be compressed into more than $NH(X)$ bits with negligible risk of information loss, as $N \rightarrow \infty$; but conversely, if they are compressed into fewer than $NH(X)$ bits it is virtually certain that information will be lost.*

This fundamental theorem specifies a theoretical limit of possibility for lossless compression. Compression at rates close to the entropy is possible (though may not be computationally feasible), and coding at better rates, without loss of information, is not.

1.2.2 Coding according to a model

In the next section, some of the coding algorithms which we review assume access to a probabilistic model, as defined in Section 1.1. In practical settings, there is almost always some discrepancy between the model used and the true data generating distribution. To highlight this, we will use $\tilde{p}_i = P(x = a_i; \theta)$ and $h(x = a_i; \theta) = \log 1/\tilde{p}_i$ for the probabilities and information content *according to the model*. The coding algorithms are deterministic, and for a fixed input sequence their message length has no dependence on the true generative distribution, but it will depend on the model distribution, and in particular on $h(x; \theta)$. In particular, we will see that, for the three model-based algorithms that we describe, we can write down a bound on compressed message length of the form

$$l(x) \leq h(x; \theta) + \dots \quad (1.7)$$

This is useful, because the maximum likelihood objective, defined in eq. (1.1), directly minimizes $h(x; \theta)$ on training data. Thus we might say that the log-likelihood is the correct objective function to use to optimize a model for lossless compression (assuming any extra terms, denoted ' \dots ', are not too significant).

We can relate the above quantity to the entropy by taking the expectation of $h(x; \theta)$ over the true generating distribution

$$\mathbb{E}(h(x; \theta)) = \sum_i p_i \log \frac{1}{\tilde{p}_i} \quad (1.8)$$

$$= H(X) + \sum_i p_i \log \frac{p_i}{\tilde{p}_i}. \quad (1.9)$$

The sum $\sum_i p_i \log(p_i/\tilde{p}_i)$ is an example of a Kullback-Leibler divergence (KL divergence). It can be shown, using Jensen's inequality, that it is non-negative, and equal to zero only when $p_i = \tilde{p}_i$ for all i , i.e. only when the model perfectly fits the true generating distribution (see MacKay, 2003 for more detail). For an in-depth overview of explicit probabilistic model-based approaches to lossless compression, see Steinruecken (2014).

1.3 Basic coding techniques

We now give an overview of algorithms for lossless compression and the compression rates which they achieve. For Huffman coding, arithmetic coding (AC) and dictionary coding we give only brief overviews, since all of these algorithms are already described in many textbooks including MacKay (2003) and Cover and Thomas (1991). The more recently invented asymmetric numeral systems (ANS) is introduced briefly here and described in detail in Chapter 2.

1.3.1 Huffman coding

Huffman coding (Huffman, 1952) is a ‘symbol code’, which means it directly maps individual symbols, from an alphabet \mathcal{A}_X , to binary words. Huffman coding assumes access to a model ensemble over symbols, i.e. a set of probabilities $\tilde{p}_1, \dots, \tilde{p}_I$. A string containing concatenated Huffman code words can be unambiguously decoded because the words which the Huffman coder outputs have the ‘prefix property’, which means that no binary code word is equal to the start of any other code word. For example, it is possible for the codewords 0 and 10 to be output by a Huffman coder, but then the codeword 010 would not be permitted, because it would be impossible to disambiguate this word from the concatenation of the two words 0 and 10.

Given an ensemble X , Huffman coding uses an efficient iterative algorithm to generate codewords with length l_i satisfying

$$l_i = \lceil h(x = a_i; \theta) \rceil < h(x = a_i; \theta) + 1, \quad (1.10)$$

where $\lceil x \rceil$ denotes the nearest integer greater than or equal to x . The average length is therefore equal to the entropy $H(X)$ when: (a) the model is equal to the true generating distribution and (b) all of the probabilities are exact integer powers of 2, i.e. when each $p_i = 2^{-n_i}$ for $n_i \in \{0, 1, \dots\}$. A Huffman code is optimal within the class of symbol codes (codes that directly map symbols to code words) assuming only condition (a); i.e. even when the average codeword length is not equal to the entropy, it is not possible to do any better than Huffman, this is a consequence of the fact that codewords lengths must be whole numbers, and the Huffman lengths are the smallest integers greater than or equal to the information content (MacKay, 2003).

Despite achieving optimal compression rates for individual symbols, Huffman coding can be highly inefficient for compressing *sequences* of symbols, particularly when the information content of each individual symbol in a sequence is close to zero, when the contribution from the '+ 1' in the upper bound in eq. (1.10) can dominate, leading to a total average message length that is much larger than the sequence's entropy. This inefficiency is addressed by stream codes, which achieve near-optimal per-symbol compression rates over sequences of symbols, rather than for single symbols.

1.3.2 Stream codes

There are two broad classes of codes which work on sequential or streaming data. The first, which we will refer to as 'model-based stream codes', are, in a sense, a direct generalization of Huffman codes to sequential data. Like Huffman codes, they require access to a model over data, but unlike Huffman codes they cater particularly to auto-regressive (sometimes referred to as 'adaptive') models over sequences of symbols. To be precise, for data drawn from a sequence X_1, \dots, X_N , they assume access to the CDFs and inverse CDFs under the conditional distributions $P(x_1), P(x_2 | x_1), \dots, P(x_N | x_1, \dots, x_{N-1})$.

The second major class of stream codes used in practice are known as 'dictionary codes'; they require no prior knowledge, or model, of the distribution of the input stream and aim to achieve acceptable (as opposed to optimal) performance for *any* input.

1.3.2.1 Arithmetic coding and asymmetric numeral systems

'Arithmetic coding' (AC; Witten et al., 1987) and the much more recent 'asymmetric numeral systems' (ANS; Duda, 2009) are both model-based stream codes,

which work by encoding a sequence of symbols one-at-a-time, and differ from each other in the order in which data may be decoded. AC is first-in-first-out (FIFO) or ‘queue-like’, so data are decoded in the *same* order in which they are encoded, whilst ANS is last-in-first-out (LIFO) or ‘stack-like’, which means that in ANS, data are decoded in the *opposite* order to that in which they were encoded. Both schemes have very similar compression performance, achieving per-symbol compression rates which are close to the information content under the model. MacKay (2003) shows that AC can achieve a message length with the worst-case bound

$$l(x) \stackrel{\text{approx}}{\leq} h(x; \theta) + 2, \quad (1.11)$$

where $h(x; \theta)$ denotes the information content of an entire sequence, which can be decomposed as

$$h(x; \theta) = \sum_{n=1}^N h(x_n | x_1, \dots, x_{n-1}; \theta). \quad (1.12)$$

The conditional information contents $h(x_n | x_1, \dots, x_{n-1}; \theta)$ are defined in the obvious way as the information content of x_n under the conditional distribution $P(x_n | x_1, \dots, x_{n-1})$. For long sequences, the effect of the ‘+2’ on the per-symbol compression rate becomes negligible, hence we can say that the per-symbol compression rate of AC is close to optimal.

The ‘approx’ above the ‘ \leq ’ symbol is there because Mackay, and other works which mention this bound, such as Witten et al. (1987) and Moffat et al. (1998), assume the implementation can use exact (i.e. unbounded precision) rational numbers. With exact arithmetic, AC encode and decode compute time scales poorly with the sequence length N . In practical settings it is important to achieve fast runtimes, and thus implementations of AC which are used in production invariably use faster, fixed precision arithmetic to approximate the exact algorithm. Another difference between practical implementations and the ideal implementation to which eq. (1.11) applies is that files are usually comprised of a whole-number of bytes. We can easily account for this and make the bound more realistic by rounding the quantity on the right-hand-side of eq. (1.11) up to the nearest multiple of 8:

$$l(x) \stackrel{\text{approx}}{\leq} \lceil h(x; \theta) + 2 \rceil_8 \leq h(x; \theta) + 10, \quad (1.13)$$

where we use $\lceil x \rceil_n$ to denote the smallest multiple of n which is greater than or equal to x . I have been unable to find analysis of the worse-case behavior of AC when using more realistic bounded-precision arithmetic, but it is likely that eq. (1.13) is a reasonably good approximation in most practical settings.

On the other hand, establishing exact worst-case bounds for practical implementations of ANS, the algorithm that we use throughout this work, is fairly straightforward, and in Chapter 2 we show that under our ANS implementation we have

$$l(x) \leq h(x) + N\epsilon + C, \quad (1.14)$$

where ϵ and C are both implementation dependent constants. A typical setup in our experiments gives $\epsilon \approx 2.2 \times 10^{-5}$ and $C = 64$. Although this bound appears worse than the AC bound in eq. (1.13), for long enough sequences the difference has little practical implication—the ‘one-off’ constant C has a negligible effect on the per-symbol compression rate, and ϵ amounts to a worst-case overhead of one bit every 2.2×10^5 symbols. The expected (as opposed to worst-case) behaviour of ANS has not been exactly characterized but empirically it is usually significantly better than the above bound (Duda, 2009).

In practice, ANS is usually slightly faster than AC, is more straightforward to implement and is also easier to generalize to a vectorized implementation (Duda, 2009; Giesen, 2014). Moreover, the LIFO property of ANS is critical for the new methods introduced in this thesis. In Chapter 2 we describe ANS in detail, and we give a full working Python implementation, which is only 50 lines in length, in Appendix A. We discuss vectorization, which we have also implemented, in Chapter 4. A high level description of AC can be found in MacKay (2003), Chapter 6, and a detailed description with a full, clear C implementation in Witten et al. (1987).

1.3.2.2 Dictionary coding

Dictionary coding, also referred to as substitution coding, works by going through a stream of symbols and replacing (‘substituting’) any sub-sequence which has already occurred in the stream with a pointer to the sub-sequence’s previous occurrence. Compression is achieved when sequences are repeated whose length exceeds that of a pointer. The family of Lempel-Ziv coders use this technique, and a variant called DEFLATE is used in `gzip`, which is perhaps the most widely used compression software in existence. Although they do not make

use of a model over data, Lempel-Ziv codes can be shown to be asymptotically optimal (Cover and Thomas, 1991). However, for practical sources and finite sequences, the lengths of encoded messages are often significantly greater than the entropy. See MacKay (2003), Section 6.4 for more detail, including examples of sources on which Lempel-Ziv coding performs poorly.

1.4 Image compression in practice

In Chapters 3 and 5 we present the BB-ANS algorithm, which extends ANS to latent variable models. The algorithm is generic, in the sense that it can be applied with a wide range of latent variable models, over any kind of data. We chose to demonstrate the method on images because deep generative models of images are well studied and reasonably straightforward to setup and train, and because image compression is extremely widely used and well studied, with good baselines to compare to.

For most image compression use cases, such as communication of images across the internet and local storage of a photo library, some loss of data is acceptable, particularly if this can be achieved without affecting the perceptual quality of the image (i.e. the appearance of the image ‘to the human eye’). JPEG (G. K. Wallace, 1991) has been the dominant lossy format for compressing photographic images since its introduction in 1992, and is extremely widely used on the web and as a storage format.

Lossless image compression is useful when it is not known in advance that degradation in an image’s quality will be acceptable to the intended recipient. When using lossy compression there is always a somewhat speculative decision that needs to be made about how much compression is appropriate. As argued in Sneyers and Wuille (2016), it is useful to have a codec that users don’t need to think about before using, and lossless codecs fall into this category. Moreover, lossless image codecs are typically not specialised to photographs (as JPEG is), and aim to achieve reasonable performance on any image that a human might wish to communicate or store.

Lossless compression is also particularly useful for saving intermediate versions during image editing, where repeatedly editing and saving in a lossy format would lead to successively more severe degradation of image quality. It is also especially desirable for archival storage, and for images which are used in scientific or medical applications; indeed in some jurisdictions it is illegal for

medical images to be stored in a lossy format (Liu et al., 2017).

1.4.1 Comparing PNG, FLIF, and our new codecs

We now discuss two existing formats, and how the approach they take differs from ours. The first is ‘portable network graphics’ (PNG), which is the most widely used lossless image format today, and is particularly popular for compressing graphics for the web. PNG development began in the mid 1990s, and the format was intended to be a more flexible replacement for GIF (another lossless format), and to be completely unencumbered by patents, which GIF, at the time, was not. The authors of the format claimed at the time that PNG compression was “among the best that can be had without losing image data and without paying patent or other licensing fees” (Roelofs, 1999).

The reference PNG implementation, `libpng`, is free software, and the PNG standard and `libpng` were developed by a small community which overlapped with the group who developed `gzip`. PNG compression has two stages, which are referred to as ‘preprocessing’ and ‘compression’. The compression stage takes the result of preprocessing and compresses it using the same algorithm (and in `libpng` the same implementation) as `gzip`. Thus the job of preprocessing is to reversibly convert the image to a more ‘`gzip` friendly’ form.

The preprocessing in PNG consists of passing one of five convolutional filters over each line of the image. The filters all subtract the values of previous pixels from the current pixel’s value, exploiting the approximate smoothness of images to try to create a result in which all pixel values are close to zero. The result of preprocessing, along with a list specifying which of the five filters was used for each row, is passed forwards to the compression stage, which uses the same DEFLATE algorithm as `gzip`, a form of dictionary coding (Roelofs, 1999). The method for selecting which filter should be used for each row is not specified by the PNG standard, but the implementation in `libpng` uses a simple heuristic, simply selecting the filter whose output is closest to zero (in the sense that the sum of the absolute values of each element in the row is smallest).

Figure 1.1 shows the importance of the preprocessing used by PNG, since PNG is 23% smaller than `gzip` in this case. Displaying a bitstream in this way is a crude but useful way to observe when a stream contains obvious redundancy, implying that there is room for improvement in the codec. A perfect codec should output bits which are indistinguishable from samples from an i.i.d. uniform

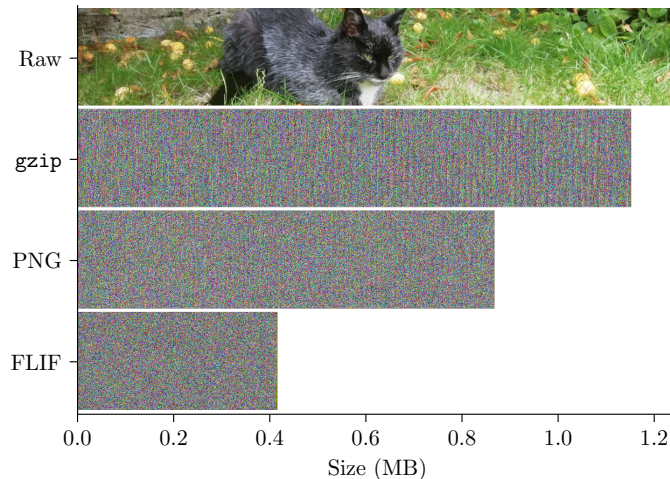


Figure 1.1: Comparing a photograph of the cat ‘Stasha’ with the compressed bit streams resulting from applying the codecs `gzip`, PNG and FLIF to the image. Vertical bands are clearly visible in the `gzip` stream, and upon close inspection there are some visible vertical bands in the PNG stream, implying some correlation between pixels and hence redundancy. No structure is visible in the FLIF stream, suggesting that FLIF may be close to optimal.

distribution (MacKay, 2003). It’s impossible to confirm that a stream is i.i.d. uniform using this method, but sometimes non-uniformities and correlation are obvious enough to be confident that a stream is *not* i.i.d. uniform.

The second format which we discuss is the ‘free lossless image format’ (FLIF), which was the last to set new state of the art performance on lossless compression benchmarks, improving on PNG by 43% on average (Sneyers and Wuille, 2016). FLIF is also free software, though it has yet to gain widespread adoption¹. It is significantly more sophisticated than PNG, although it also uses a two stage preprocessing and compression pipeline. In the FLIF documentation, preprocessing is referred to as ‘prediction’, but for self-consistency we will continue to use ‘preprocessing’.

Preprocessing in FLIF is actually slightly simpler than in PNG, in that the same filter is used on every line of the image, instead of switching between five filters. The filter works by subtracting the median of T , L and $T + L - TL$ from each pixel, where T is the value of the pixel above, and L is the value of the pixel to the left of the current pixel. The ‘compression’ stage uses an algorithm called ‘meta-adaptive near-zero integer arithmetic coding’ (MANIAC). MANIAC uses arithmetic coding with an adaptive model, which learns as it processes an image.

¹In fact FLIF has recently been subsumed by the JPEG XL standard, still under development at time of writing, see Alakuijala et al., 2019.

Which mixture component is selected for each pixel is determined by the values of the pixels T and L , see Sneyers and Wuille (2016) for more detail.

1.4.1.1 A paradigm shift

The approach we take to lossless compression is fundamentally different to PNG, FLIF, and other existing codecs. We explicitly de-couple compression and modelling, spending a lot of engineering effort and compute time to develop an accurate image model, *before* doing compression with the model. We can train our model on a dataset containing a representative sample of the images that the compression algorithm will be used for. For the full-size colour image compression experiments in Chapter 5, we use the ImageNet dataset, which is reasonably representative of the images that are communicated over the internet. Having ‘seen’ a huge number of images before attempting any compression, our method could be said, in the language of Bayesian statistics, to have a ‘strong prior’. Before compressing a test image, it may have very high level structural knowledge about the kinds of objects that tend to appear in images. Given an image whose top few rows look like the top of a cat’s head, the model may be able to infer that the rest of the image is likely to contain the rest of the cat. PNG, FLIF and other existing codecs could only be said to have weak knowledge about images—they know that pixels are locally correlated and that patterns tend to be repeated within individual images. They certainly would not be able to predict the existence of pixels representing a cat’s body given only the top of the cat’s head.

This approach represents a paradigm shift, and in Chapters 3 and 5 we demonstrate that this may lead to improved compression rates. Currently, our implementations are significantly slower than existing codecs, and the BB-ANS method only performs optimally when processing *batches* of images, but we are confident that these issues can be overcome, in future work which builds on the foundation we have helped to lay.

Chapter 2

Introduction to asymmetric numeral systems

We are interested in algorithms for lossless compression of sequential data. Arithmetic coding (AC) and the range variant of asymmetric numeral systems (sometimes abbreviated to rANS, we simply use ANS) are examples of such algorithms. Just like arithmetic coding, ANS is close to optimal in terms of compression rate (Witten et al., 1987; Duda, 2009). The key difference between ANS and AC is in the order in which data are *decoded*: in ANS, compression is last-in-first-out (LIFO), or ‘stack-like’, while in AC it is first-in-first-out (FIFO), or ‘queue-like’. The stack-like nature of ANS is critical for the algorithm we present in Chapter 3. We recommend MacKay (2003) Chapter 4-6 for background on source coding and arithmetic coding in particular. This chapter contains pseudocode which could be converted to a working ANS implementation without too much difficulty. We also provide a 50 line Python implementation, with example usage, in Appendix A.

ANS comprises two basic functions, which we denote `push` and `pop`, for encoding and decoding, respectively (the names refer to the analogous stack operations). The `push` function accepts some pre-compressed information m (short for ‘message’), and a symbol x to be compressed, and returns a new compressed message, m' . Thus it has the signature

$$\text{push}: (m, x) \mapsto m'. \quad (2.1)$$

The new compressed message, m' , contains precisely the same information as the pair (m, x) , and therefore `push` can be inverted to form a decoder mapping.

The decoder, `pop`, maps from m' back to m, x :

$$\text{pop}: m' \mapsto (m, x). \quad (2.2)$$

Because the functions `push` and `pop` are inverse to one another, we have $\text{push}(\text{pop}(m)) = m$ and $\text{pop}(\text{push}(m, x)) = (m, x)$.

2.1 Specifying the problem which ANS solves

In this section we first define some notation, then describe the problem which ANS solves in more detail and sketch the high level approach to solving it. In the following we use ‘log’ as shorthand for the base 2 logarithm, usually denoted ‘ \log_2 ’.

The functions `push` and `pop` will both require access to the probability distribution from which symbols are drawn (or an approximation thereof). To describe distributions we use notation similar to MacKay (2003):

Definition 4. A *quantized ensemble* X with precision r is a triple $(x, \mathcal{A}_X, \mathcal{P}_X)$ where the *outcome* x is the value of a random variable, taking on one of a set of possible values $\mathcal{A}_X = \{a_1, \dots, a_I\}$, and $\mathcal{P}_X = \{p_1, \dots, p_I\}$ are the *integer-valued* probability weights with each $p_i \in \{1, \dots, 2^r\}$, each $P(x = a_i) = p_i/2^r$ and therefore $\sum_{i=1}^I p_i = 2^r$.

Note that this definition differs from Definition 1 from the last chapter, in that the probabilities are assumed to be *quantized* to some precision r (i.e. representable by fractions $p_i/2^r$), and we assume that none of the a_i have zero probability. Having probabilities in this form is necessary for the arithmetic operations involved in ANS (as well as AC). Note that if we use a high enough r then we can specify probabilities that are not close to zero with a precision similar to that of typical floating point—32-bit floating point numbers for example contain 23 ‘fraction’ bits, and thus would have roughly the same precision as our representation with $r = 23$. Symbols with very small probabilities are a potential failure mode of both ANS and AC. These must either be rounded up to the smallest nonzero quantized value, or all symbols with small probabilities can effectively be split into two symbols: the first, a global ‘escape’ symbol, which has mass equal to the sum of the small probabilities, and the second, the symbol itself, with mass equal to the conditional given that the symbol is in the escape group (this conditional will be larger than the symbol’s original mass).

An escape mechanism is described in more detail in Moffat et al. (1998).

Given a sequence of quantized ensembles X_1, \dots, X_N , we seek an algorithm which can encode any outcome x_1, \dots, x_N in a binary message whose length is close to $h(x_1, \dots, x_N) = \log 1/P(x_1, \dots, x_N)$. According to Shannon's source coding theorem it is not possible to losslessly encode data in a message with expected length less than $\mathbb{E}[h(x)]$, thus we are looking for an encoding which is close to optimal in expectation (Shannon, 1948). Note that the joint information content of the sequence can be decomposed:

$$h(x_1, \dots, x_N) = \log \frac{1}{P(x_1, \dots, x_N)} \quad (2.3)$$

$$= \sum_n \log \frac{1}{P(x_n | x_1, \dots, x_{n-1})} \quad (2.4)$$

$$= \sum_n h(x_n | x_1, \dots, x_{n-1}). \quad (2.5)$$

Because it simplifies the presentation significantly, we focus first on the ANS *decoder*, the reverse mapping which maps from a compressed binary message to the sequence x_1, \dots, x_N . This will be formed of a sequence of N **pop** operations; starting with a message m_0 we define

$$m_n, x_n = \mathbf{pop}(m_{n-1}) \quad \text{for } n = 1, \dots, N \quad (2.6)$$

where each **pop** uses the conditional distribution $X_n | X_1, \dots, X_{n-1}$. We will show that the message resulting from each **pop**, m_n , is effectively shorter than m_{n-1} by no more than $h(x_n | x_1, \dots, x_{n-1}) + \epsilon$ bits, where ϵ is a small constant which we specify below, and therefore the difference in length between m_0 and m_N is no more than $h(x_1, \dots, x_N) + N\epsilon$, by eqs. (2.3) to (2.5).

We will also show that **pop** is a bijection whose inverse, **push**, is straightforward to compute, and therefore an encoding procedure can easily be defined by starting with a very short base message and adding data sequentially using **push**. Our guarantee about the effect of **pop** on message length translates directly to a guarantee about the effect of **push**, in that the increase in message length due to the sequence of **push** operations is less than $h(x_1, \dots, x_N) + N\epsilon$.

2.2 Asymmetric numeral systems

Having set out the problem which ANS solves and given a high level overview of the solution in Section 2.1, we now go into more detail, firstly discussing the

data structure we use for m , then the pop function and finally the computation of its inverse, `push`.

2.2.1 The structure of the message

We use a pair $m = (s, t)$ as the data structure for the message m . The element s is an unsigned integer with precision r_s (i.e. $s \in \{0, 1, \dots, 2^{r_s} - 1\}$, so that s can be expressed as a binary number with r_s bits). The element t is a stack of unsigned integers of some fixed precision r_t where $r_t < r_s$. This stack has its own push and pop operations, which we denote `stack_push` and `stack_pop` respectively. See fig. 2.1 for a diagram of s and t . We need s to be large enough to ensure that our decoding is accurate, and so we also impose the constraint

$$s \geq 2^{r_s - r_t}, \quad (2.7)$$

more detail on how and why we do this is given below. In the demo implementation we use $r_s = 64$ and $r_t = 32$.

Note that a message can be flattened into a string of bits by concatenating s and the elements of t . The length of this string is

$$l(m) := r_s + r_t |t|, \quad (2.8)$$

where $|t|$ is the number of elements in the stack t . We refer to this quantity as the ‘length’ of m . We also define the useful quantity

$$l^*(m) := \log s + r_t |t|, \quad (2.9)$$

which we refer to as the ‘effective length’ of m . Note that the constraint in eq. (2.7) and the fact that $s < 2^{r_s}$ imply that

$$l(m) - r_t \leq l^*(m) < l(m). \quad (2.10)$$

Intuitively l^* can be thought of as a precise measure of the size of m , whereas l , which is integer valued, is a more crude measure. Clearly l is ultimately the measure that we care most about, since it tells us the size of a binary encoding of m , and we use l^* to prove bounds on l .

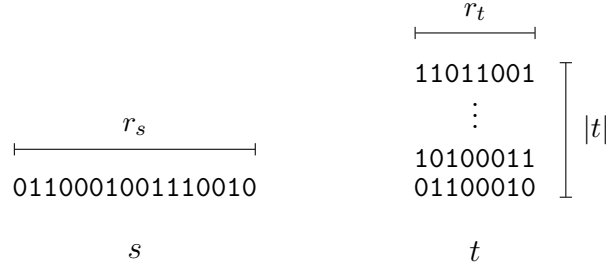


Figure 2.1: The two components of a message: the unsigned integer s (with $r_s = 16$) and the stack of unsigned integers t (with $r_t = 8$). The integers are represented here in base 2 (binary).

2.2.2 Constructing the pop operation

To avoid notational clutter, we begin by describing the `pop` operation for a single quantized ensemble $X = (x, \mathcal{A}_X, \mathcal{P}_X)$ with precision r , before applying `pop` to a sequence in Section 2.2.3. Our strategy for performing a decode with `pop` will be firstly to extract a symbol from s . We do this using a bijective function $d: \mathbb{N} \rightarrow \mathbb{N} \times \mathcal{A}$, which takes an integer s as input and returns a pair (s', x) , where s' is an integer and x is a symbol. Thus `pop` begins

```
def pop(m):
    s, t := m
    s', x := d(s)
```

We design the function d so that if $s \geq 2^{r_s - r_t}$, then

$$\log s - \log s' \leq h(x) + \epsilon \quad (2.11)$$

where

$$\epsilon := \log \frac{1}{1 - 2^{-(r_s - r_t - r)}}. \quad (2.12)$$

We give details of d and prove eq. (2.11) below. Note that when the term $2^{-(r_s - r_t - r)}$ is small, the following approximation is accurate:

$$\epsilon \approx \frac{2^{-(r_s - r_t - r)}}{\ln 2}, \quad (2.13)$$

and thus ϵ itself is small. We typically use $r_s = 64$, $r_t = 32$, and $r = 16$, which gives $\epsilon = \log 1/(1 - 2^{-16}) \approx 2.2 \times 10^{-5}$.

After extracting a symbol using d , we check whether s' is below $2^{r_s - r_t}$, and if it is we `stack_pop` integers from t and move their contents into the lower order bits of s' . We refer to this as ‘renormalization’. Having done this, we return the new message and the symbol x . The full definition of `pop` is thus

```

def pop(m):
    s, t := m
    s', x := d(s)
    s, t := renorm(s', t)
    return (s, t), x

```

Renormalization is necessary to ensure that the value of s returned by `pop` satisfies $s \geq 2^{r_s - r_t}$ and is therefore large enough that eq. (2.11) holds at the start of any future `pop` operation. The `renorm` function has a while loop, which pushes elements from t into the lower order bits of s until s is full to capacity. To be precise:

```

def renorm(s, t):
    # while s has space for another element from t
    while s < 2^{r_s - r_t}:
        # pop an element t_top from t
        t, t_top := stack_pop(t)
        # and push t_top into the lower bits of s
        s := 2^{r_t} * s + t_top
    return s, t

```

The condition $s < 2^{r_s - r_t}$ guarantees that $2^{r_t} \cdot s + t_{\text{top}} < 2^{r_s}$, and thus there can be no loss of information resulting from overflow. We also have

$$\log(2^{r_t} \cdot s + t_{\text{top}}) \geq r_t + \log s \quad (2.14)$$

since $t_{\text{top}} \geq 0$. Applying this inequality repeatedly, once for each iteration of the while loop in `renorm`, we have

$$\log s \geq \log s' + r_t \cdot [\# \text{ elements popped from } t], \quad (2.15)$$

where $s, t = \text{renorm}(s', t)$ as in the definition of `pop`.

Combining eq. (2.11) and eq. (2.15) gives us

$$l^*(m) - l^*(m') \leq h(x) + \epsilon, \quad (2.16)$$

where $(m', x) = \text{pop}(m)$, using the definition of l^* . That is, the reduction in the effective message length resulting from `pop` is close to $h(x)$.

2.2.3 Popping in sequence

We now apply `pop` to the setup described in Section 2.1, performing a sequence of `pop` operations to decode a sequence of data. We suppose that we are given

some initial message m_0 .

For $n = 1 \dots N$, we let $m_n, x_n = \text{pop}(m_{n-1})$ as in Section 2.1, where each `pop` uses the corresponding distribution $X_n | X_1, \dots, X_{n-1}$. Applying eq. (2.16) to each of the N `pop` operations, we have:

$$l^*(m_0) - l^*(m_N) = \sum_{n=1}^N [l^*(m_{n-1}) - l^*(m_n)] \quad (2.17)$$

$$\leq \sum_{n=1}^N [h(x_n | x_1, \dots, x_{n-1}) + \epsilon] \quad (2.18)$$

$$\leq h(x_1, \dots, x_N) + N\epsilon. \quad (2.19)$$

This result tells us about the reduction in message length from `pop` but also, conversely, about the length of a message *constructed* using `push`. We can actually initialize an encoding procedure by *choosing* m_N , and then performing a sequence of `push` operations. Since our ultimate goal when encoding is to minimize the encoded message length m_0 we choose the setting of m_N which minimizes $l^*(m_N)$, which is $m_N = (s_N, t_N)$ where $s_N = 2^{r_s - r_t}$ and t_N is an empty stack. That gives $l^*(m_N) = r_s - r_t$ and therefore, by eq. (2.19),

$$l^*(m_0) \leq h(x_1, \dots, x_N) + N\epsilon + r_s - r_t. \quad (2.20)$$

Combining that with eq. (2.10) gives an expression for the actual length of the flattened binary message resulting from m_0 :

$$l(m_0) \leq h(x_1, \dots, x_N) + N\epsilon + r_t. \quad (2.21)$$

It now remains for us to describe the function d and show that it satisfies eq. (2.11), as well as showing how to invert `pop` to form the encoding function `push`.

2.2.4 The function d

The function $d: \mathbb{N} \rightarrow \mathbb{N} \times \mathcal{A}$ must be a bijection, and we aim for d to satisfy eq. (2.11), and thus $P(x) \approx \frac{s'}{s}$. Achieving this is actually fairly straightforward. One way to define a bijection $d: s \mapsto (s', x)$ is to start with a mapping $\tilde{d}: s \mapsto x$, with the property that none of the preimages $\tilde{d}^{-1}(x) := \{n \in \mathbb{N} : \tilde{d}(n) = x\}$ are finite for $x \in \mathcal{A}$. Then let s' be the index of s within the (ordered) set $\tilde{d}^{-1}(x)$, with indices starting at 0. Equivalently, s' is the number of integers n with

$0 \leq n < s$ and $d(n) = x$.

With this setup, the ratio

$$\frac{s'}{s} = \frac{|\{n \in \mathbb{N} : n < s, d(n) = x\}|}{s} \quad (2.22)$$

is the density of numbers which decode to x , within all the natural numbers less than s . For large s we can ensure that this ratio is close to $P(x)$ by setting \tilde{d} such that numbers which decode to a symbol x are distributed *within the natural numbers* with density close to $P(x)$.

To do this, we partition \mathbb{N} into finite ranges of equal length, and treat each range as a model for the interval $[0, 1]$, with sub-intervals within $[0, 1]$ corresponding to each symbol, and the width of each sub-interval being equal to the corresponding symbol's probability (see fig. 2.2). To be precise, the mapping \tilde{d} can then be expressed as a composition $\tilde{d} = \tilde{d}_2 \circ \tilde{d}_1$, where \tilde{d}_1 does the partitioning described above, and \tilde{d}_2 assigns numbers within each partition to symbols (sub-intervals). So

$$\tilde{d}_1(s) := s \bmod 2^r. \quad (2.23)$$

Using the shorthand $\bar{s} := \tilde{d}_1(s)$, and defining

$$c_j := \begin{cases} 0 & \text{if } j = 1 \\ \sum_{k=1}^{j-1} p_k & \text{if } j = 2, \dots, I \end{cases} \quad (2.24)$$

as the (quantized) cumulative probability of symbol a_{j-1} ,

$$\tilde{d}_2(\bar{s}) := a_i \text{ where } i := \max\{j : c_j \leq \bar{s}\}. \quad (2.25)$$

That is, $\tilde{d}_2(\bar{s})$ selects the symbol whose sub-interval contains \bar{s} . Figure 2.2 illustrates this mapping, with a particular probability distribution, for the range $s = 64, \dots, 71$.

2.2.5 Computing s'

The number s' was defined above as “the index of s within the (ordered) set $\tilde{d}^{-1}(x)$, with indices starting at 0”. We now derive an expression for s' in terms of s , p_i and c_i , where $i = \max\{j : c_j \leq \bar{s}\}$ (as above), and we prove eq. (2.11).

Our expression for s' is a sum of two terms. The first term counts the

s	64	65	66	67	68	69	70	71	
$s \bmod 2^r$	0	1	2	3	4	5	6	7	
x	a	b			c		d		
	0								1

Figure 2.2: Showing the correspondence between s , $s \bmod 2^r$ and the symbol x . The interval $[0, 1] \subset \mathbb{R}$ is modelled by the set of integers $\{0, 1, \dots, 2^r - 1\}$. In this case $r = 3$ and the probabilities of each symbol are $P(\mathbf{a}) = 1/8$, $P(\mathbf{b}) = 2/8$, $P(\mathbf{c}) = 3/8$ and $P(\mathbf{d}) = 2/8$.

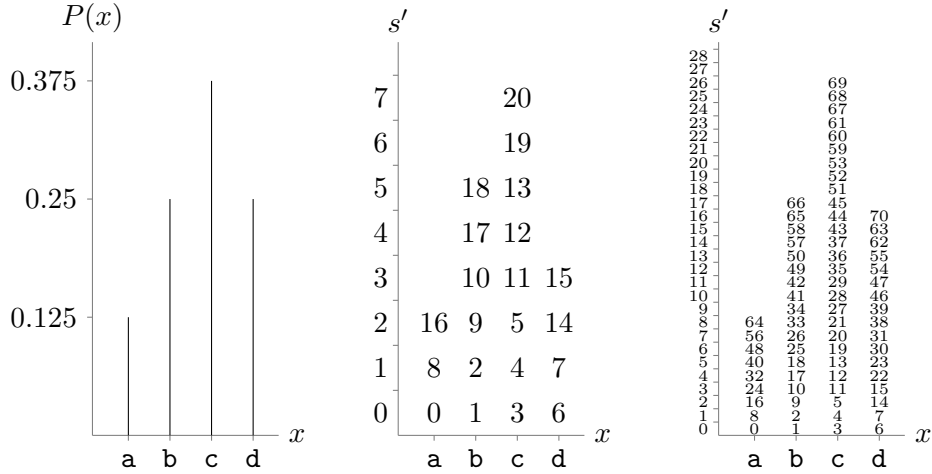


Figure 2.3: Showing the pmf of a distribution over symbols (left) and a visualization of the mapping d (middle and right). In the middle and right figures, numbers less than or equal to s_{\max} are plotted, for $s_{\max} = 20$ and $s_{\max} = 70$. The position of each number s plotted is set to the coordinates (x, s') , where $s', x = d(s)$. The heights of the bars are thus determined by the ratio s'/s from eq. (2.22), and can be seen to approach the heights of the lines in the histogram on the left (that is, to approach $P(x)$) as the density of numbers increases.

entire intervals, corresponding to the selected symbol a_i , which are below s . The size of each interval is p_i and the number of intervals is $s \div 2^r$, thus the first term is $p_i \cdot (s \div 2^r)$, where \div denotes *integer* division, discarding any remainder. The second term counts our position within the current interval, which is $\bar{s} - c_i \equiv s \bmod 2^r - c_i$. Thus

$$s' = p_i \cdot (s \div 2^r) + s \bmod 2^r - c_i. \tag{2.26}$$

This expression is straightforward to compute. Moreover from this expression it is straightforward to prove eq. (2.11). Firstly, taking the log of both sides of eq. (2.26) and using the fact that $s \bmod 2^r - c_i \geq 0$ gives

$$\log s' \geq \log(p_i \cdot (s \div 2^r)). \tag{2.27}$$

then by the definition of \div , we have $s \div 2^r > \frac{s}{2^r} - 1$, and thus

$$\log s' \geq \log \left(p_i \left(\frac{s}{2^r} - 1 \right) \right) \quad (2.28)$$

$$\geq \log s - h(x) + \log \left(1 - \frac{2^r}{s} \right) \quad (2.29)$$

$$\geq \log s - h(x) - \epsilon, \quad (2.30)$$

as required, using the fact that $P(x) = \frac{p_i}{2^r}$ and $s \geq 2^{r_s - r_t}$.

By choosing $r_s - r_t$ to be reasonably large (it is equal to 32 in our implementation), we ensure that $\frac{s'}{s}$ is very close to $P(x)$. This behaviour can be seen visually in fig. 2.3, which shows the improvement in the approximation for larger s .

2.2.6 Pseudocode for d

We now have everything we need to write down a procedure to compute d . We assume access to a function $f_X: \bar{s} \mapsto (a_i, c_i, p_i)$, where i is defined above. This function clearly depends on the distribution of X , and its computational complexity is equivalent to that of computing the CDF and inverse CDF for X . For many common distributions, the CDF and inverse CDF have straightforward closed form expressions, which don't require an explicit sum over i .

We compute d as follows:

```
def d(s):
     $\bar{s} := s \bmod 2^r$ 
     $x, c, p := f_X(\bar{s})$ 
     $s' := p \cdot (s \div 2^r) + \bar{s} - c$ 
    return  $s', x$ 
```

2.2.7 Inverting the decoder

Having described a decoding process which appears not to throw away any information, we now derive the inverse process, `push`, and show that it is computationally straightforward.

The `push` function has access to the symbol x as one of its inputs, and must do two things. Firstly it must `stack_push` the correct number of elements to t from the lower bits of s . Then it must reverse the effect of d on s , returning a value of s identical to that before `pop` was applied.

Thus, on a high level, the inverse of the function `pop` can be expressed as

```
def push(m, x):
     $s, t := m$ 
     $p, c := g_X(x)$ 
```

```

s', t := renorm_inverse(s, t; p)
s := d^{-1}(s'; p, c)
return s, t

```

where $g_X : x \mapsto (p_i, c_i)$ with i as above. The function g_X is similar to f_X in that it is analogous to computing the quantized CDF and mass function $x \mapsto p_i$. The function d^{-1} is really a pseudo-inverse of d ; it is the inverse of $s \mapsto d(s, x)$, holding x fixed.

As mentioned above, `renorm_inverse` must `stack_push` the correct amount of data from the lower order bits of s into t . A necessary condition which the output of `renorm_inverse` must satisfy is

$$2^{r_s - r_t} \leq d^{-1}(s'; p, c) < 2^{r_s}. \quad (2.31)$$

This is because the output of `push` must be a valid message, as described in Section 2.2.1, just as the output of `pop` must be.

The expression for s' in eq. (2.26) is straightforward to invert, yielding a formula for d^{-1} :

$$d^{-1}(s'; p, c) = 2^r \cdot (s' \div p) + s' \bmod p + c. \quad (2.32)$$

We can substitute this into eq. (2.31) and simplify:

$$2^{r_s - r_t} \leq 2^r \cdot (s' \div p) + s' \bmod p + c < 2^{r_s} \quad (2.33)$$

$$\iff 2^{r_s - r_t} \leq 2^r \cdot (s' \div p) < 2^{r_s} \quad (2.34)$$

$$\iff p \cdot 2^{r_s - r_t - r} \leq s' < p \cdot 2^{r_s - r}. \quad (2.35)$$

So `renorm_inverse` should move data from the lower order bits of s' into t (decreasing s') until eq. (2.35) is satisfied. To be specific:

```

def renorm_inverse(s', t; p):
  while s' >= p * 2^{r_s - r}:
    t := stack_push(t, s' mod 2^{r_t})
    s' := s' div 2^{r_t}
  return s', t

```

Although, as mentioned above, eq. (2.35) is a *necessary* condition which s' must satisfy, it isn't immediately clear that it's sufficient. Is it possible that we need to continue the while loop in `renorm_inverse` past the first time that $s' < p \cdot 2^{r_s - r}$? In fact this can't be the case, because $s' \div 2^{r_t}$ decreases s' by a

factor of at least 2^{r^t} , and thus as we iterate the loop above we will land in the interval specified by eq. (2.35) at most once. This guarantees that the s that we recover from `renorm_inverse` is the correct one.

2.3 Further reading

Since its invention by Duda (2009), ANS appears not to have gained widespread attention in academic literature, despite being used in various state of the art compression systems. At the time of writing, a search on Google Scholar for the string “asymmetric numeral systems” yields 148 results. For comparison, a search for “arithmetic coding”, yields ‘about 44,000’ results. As far as I’m aware, ANS has appeared in only one textbook, with a practical, rather than mathematical, presentation (McAnlis and Haecky, 2016).

However, for those wanting to learn more there is a huge amount of material on different variants of ANS in Duda (2009) and Duda et al. (2015). Extensions to latent variable models are, of course, described in Chapters 3 and 5, as well as in Townsend et al. (2019), F. H. Kingma et al. (2019), and Townsend et al. (2020). A parallelized implementation based on SIMD instructions was first presented in Giesen (2014) and is described in Chapter 4. Finally, a version which performs simultaneous encryption and compression is described in Duda and Niemiec (2016).

Duda maintains a list of ANS implementations at <https://encode.su/threads/2078-List-of-Asymmetric-Numeral-Systems-implementations>.

Chapter 3

Bits back coding with asymmetric numeral systems

This chapter is based on the paper ‘Practical lossless compression with latent variables using bits back coding’ (Townsend et al., 2019), which was presented at the International Conference of Learning Representations (ICLR). The paper was co-authored with Tom Bird, who assisted with writing the paper and setting up the experiments. We propose a new coding method, extending ANS to latent variable models. Like ANS itself, our new method achieves a near-optimal rate when a sequence (or ‘batch’) of data are compressed, and may suffer a significant, but bounded, overhead at the beginning of the encoding process (we discuss methods for minimizing this overhead in Section 5.1.2). At the end of this chapter we demonstrate the method by compressing the MNIST test set using a VAE, outperforming all other codecs benchmarked, in terms of compression rate. Code for reproducing the results in this chapter can be found at github.com/bits-back/bits-back.

The lossless compression algorithms mentioned in Chapters 1 and 2, namely Huffman coding, arithmetic coding (AC) and asymmetric numeral systems (ANS), do not naturally cater for models with latent variables. However, there is a method, known as ‘bits back coding’ (C. S. Wallace, 1990; Hinton and van Camp, 1993), first introduced as a thought experiment, but later implemented in Frey and Hinton (1996) and Frey (1997), which may be used to extend those algorithms to such models.

Although bits back coding was implemented in restricted cases by Frey (1997), prior to our work there was no known implementation for modern neural net-based models or high dimensional data; Frey’s implementation was demonstrated on 8×8 binary images. There is, in fact, an awkward incompatibility

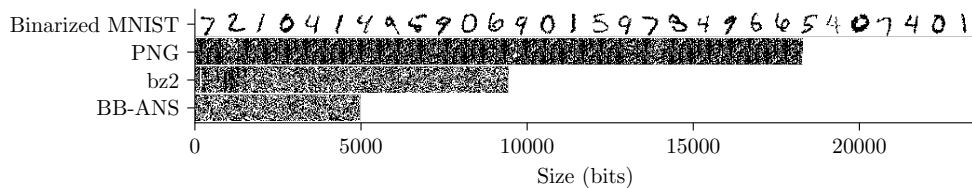


Figure 3.1: Visual comparison of 30 binarized MNIST images with bit stream outputs from running lossless compression algorithms PNG, bz2 and BB-ANS on the images.

between bits back and the arithmetic coding scheme with which it was implemented in Frey (1997). In this chapter we describe this issue and present a clean solution—a scheme that instead implements bits back coding using ANS. We term this new coding scheme ‘Bits Back with ANS’ (BB-ANS).

In Section 3.3.2 we demonstrate the efficacy of BB-ANS by losslessly compressing the MNIST dataset with a variational auto-encoder (VAE; D. P. Kingma and Welling, 2014), a deep latent variable model with continuous latent variables which we described in Section 1.1.1. BB-ANS with a VAE outperforms generic compression algorithms for both binarized and raw MNIST, even with a very simple model architecture. In Chapter 5 we demonstrate that the method scales well to a larger model, using it to compress full-size colour photographs.

3.1 Bits back coding

It is well known that an arithmetic *decoder* can be used to map a randomised binary message to a sample from the distribution used as a model for the coder. In the case of ANS, the fact that running the decoder, parametrized by some mass function P , on a random message will generate a sample from P , is a straightforward consequence of the fact, stated and proved in Section 2.2, that numbers which decode to x appear in the natural numbers with density close to $P(x)$. This is a fundamental property of ANS which is necessary for optimal compression.

In this section we describe bits back coding, a method which uses the aforementioned sampling capability of a lossless codec for compression of data using a latent variable model. We first present the form in which bits-back coding has appeared in previous works, then we present our own novel approach.

3.1.1 Bits back without ANS

Suppose that a sender wishes to communicate a symbol x to a receiver, and that both sender and receiver have access to a generative model with a latent

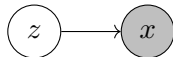


Figure 3.2: Graphical model with latent variable z and observed variable x .

variable, z . For now we take z to be discrete; in Section 3.2.2 we show how to apply the method to continuous latents.

Suppose that the CDFs/inverse CDFs required for coding with ANS under the marginal distribution $P(x)$ are intractable to compute, as is typically the case in deep latent variable models. Bits back is an idea that can be used to encode and decode x assuming only the ability to encode/decode using the forward probabilities $P(z)$ and $P(x|z)$, and a (possibly approximate) posterior $Q(z|x)$.

We must assume that, as well as the sample x , the sender already has a separate compressed message to communicate. The sender can run the approximate posterior *decode* on the extra message to generate a sample from $Q(z|x)$. Then they can encode the latent sample according to $P(z)$ and the symbol x according to $P(x|z)$. The receiver then does the inverse to recover the latent sample and the symbol. The extra message can also be precisely recovered by the receiver by *encoding* the latent sample according to $Q(z|x)$.

We can write down the expected increase in message length (over the length of the extra message at the start):

$$L(Q) = \mathbb{E}_{Q(z|x)}[-\log P(z) - \log P(x|z) + \log Q(z|x)] \quad (3.1)$$

$$= -\mathbb{E}_{Q(z|x)} \log \frac{P(x,z)}{Q(z|x)}. \quad (3.2)$$

This is the negative of the evidence lower bound (ELBO), which was defined in eq. (1.3).

A great deal of recent research has focused on inference and learning with approximate posteriors, using the ELBO as an objective function. Because of the above equivalence, methods which maximize the ELBO for a model are implicitly minimizing the message length achievable by bits back coding with that model. This suggests that we may be able to draw on this plethora of existing methods when learning a model for use in compression applications, safe in the knowledge that the objective function they are maximizing is the negative expected message length.

3.1.2 Chaining bits back coding

If we wish to encode a *sequence* of data points, and do not have any extra information to communicate, we may be able to accept a one-time overhead for coding the first element at a rate worse than the negative ELBO. Maybe we have a fallback codec which doesn't require an existing message, which we can use for the first element; if not, we can generate the latent for the first element of the sequence in any way we like, including by sampling it from a pseudo-random number generator, and then encode it using the prior and likelihood. The resulting compressed message can then be used as the extra message for bits-back coding of the second data point, the encoded second data point as the extra message for the third, and so on. This daisy-chain-like scheme was first described by Frey (1997), and was called 'bits-back with feedback'. We refer to it simply as 'chaining'.

As Frey (1997) notes, the above method cannot be implemented directly using AC, because in order for it to work it is necessary to decode data in the opposite order to that in which they were encoded. Frey gets around this by implementing what amounts to a stack-like wrapper around AC, for which it is necessary to terminate AC encoding after each chaining step, effectively using AC like a symbol code. This incurs a cost both in code complexity and, importantly, in compression rate. The cost in compression rate is due to the fact that terminating AC incurs a cost of up to two bits (see eq. 1.11). As Frey notes, any symbol code will do for chaining, and in situations where x and z are actually each comprised of individual symbols (in most situations we study, they are in fact vectors), Huffman coding would be the optimal choice, but this still incurs a compression rate overhead for each chaining step.

3.1.3 Chaining bits back coding with ANS

The central insight of this chapter is the observation that the chaining described in the previous section can be implemented straightforwardly with ANS with zero compression rate overhead per iteration. This is because of the fact that ANS is stack-like by nature, which resolves the problems that occur if one tries to implement bits back chaining with AC, which is queue-like. We now describe this novel method, which we refer to as 'Bits Back with ANS' (BB-ANS).

We can visualize the stack-like state of an ANS coder as

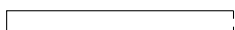
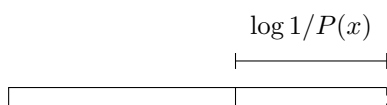


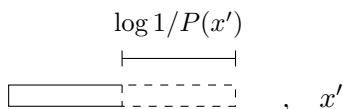
Table 3.1: The shorthand we use for ANS encoding and decoding operations, based on left and right pointing arrows. The operations can be translated mechanically to or from the pseudocode shown in the right-hand column.

Notation	Meaning	Pseudocode
$x \rightarrow P(\cdot)$	“Encode x using P ”	<code>message = push.P(message, x)</code>
$x \leftarrow P(\cdot)$	“Decode x using P ”	<code>message, x = pop.P(message)</code>

where the dashed line on the right symbolizes the encoding/decoding end or ‘top’ of the stack. When we encode a symbol x onto the message stack we effectively add it to the end, resulting in a ‘longer’ state







and when we decode (or equivalently, sample) a symbol x' from the stack we remove it from the same end, resulting in a ‘shorter’ state, plus the symbol that we decoded.



For compactness, we use define a shorthand notation for encoding and decoding operations: $x \rightarrow P(\cdot)$ for encoding (pushing) x onto the stack using the distribution P , and $x \leftarrow P(\cdot)$ for decoding (popping). The notation is summarized in Table 3.1. Table 3.2 shows the states of the message as the sender encodes a sample, using our bits back with ANS algorithm, starting with an existing message, labelled ‘extra information’, as well as the sample x to be encoded. The operations are performed starting at the top of the table and working downwards.

This process is clearly invertible, by reversing the order of operation and replacing encodes with decodes and sampling with encoding. Furthermore it can be repeated; the ANS message at the end of encoding is still an ANS message, and therefore can be readily used as the extra information for encoding the next symbol. The algorithm is compatible with any model whose prior, likelihood and (approximate) posterior can be encoded and decoded with ANS, i.e. it is necessary and sufficient to be able to compute conditional CDFs and inverse CDFs under those distributions. A simple Python implementation of both the encoder and decoder of BB-ANS is shown in fig. 3.3.

Table 3.2: Visualizing the process by which a sender encodes (pushes) a symbol x onto an ANS message stack using Bits Back with ANS. The process starts with an existing ANS message, labelled ‘extra information’, and the operations in the right hand column are performed starting at the top of the table and working downwards. The ‘Variables’ column shows variables which are known to the sender before each operation is performed.

BB-ANS message	Variables	Operation
Extra information 	x	
$\log 1/Q(z x)$ 	x, z	$z \leftarrow Q(\cdot x)$
$\log 1/P(x z)$ 	z	$x \rightarrow P(\cdot z)$
$\log 1/P(z)$ 		$z \rightarrow P(\cdot)$

3.2 Issues affecting the efficiency of BB-ANS

A number of factors can affect the efficiency of compression with BB-ANS, and mean that in practice, the coding rate will never be exactly equal to the ELBO. For any algorithm based on AC/ANS, the fact that all probabilities have to be approximated at finite precision has some detrimental effect. When encoding a batch of only a small number of i.i.d. samples, with no ‘extra information’ to communicate, the inefficiency of encoding the first data point may be significant. In the worst case, that of a batch with only one data point, the message length will be equal to the log joint, $\log P(x, z)$. Note that optimization of this is equivalent to maximum a posteriori (MAP) estimation. However, for a batch containing multiple images, this effect is amortized. Figure 3.1 shows an example with 30 samples, where BB-ANS performs well.

Below we discuss two other issues which are specific to BB-ANS. We investigate the magnitude of these effects experimentally in Section 3.3.2. We find that when compressing the MNIST test set, they do not significantly affect the compression rate, which is typically close to the negative ELBO in our experiments.

```

def push(message, x):
    # (1) Sample  $z$  according to  $Q(z|x)$ 
    #       Decreases message length by  $\log 1/Q(z|x)$ 
    message, z = posterior_pop(x)(message)

    # (2) Encode  $x$  according to the likelihood  $P(x|z)$ 
    #       Increases message length by  $\log 1/P(x|z)$ 
    message = likelihood_append(z)(message, x)

    # (3) Encode  $z$  according to the prior  $P(z)$ 
    #       Increases message length by  $\log 1/P(z)$ 
    message = prior_append(message, z)

    return message

def pop(message):
    # (3 inverse) Decode  $z$  according to  $P(z)$ 
    message, z = prior_pop(message)

    # (2 inverse) Decode  $x$  according to  $P(x|z)$ 
    message, x = likelihood_pop(z)(message)

    # (1 inverse) Decode  $z$  according to  $Q(z|x)$ 
    message = posterior_append(x)(message, z)

    return message, x

```

Figure 3.3: Python implementation of BB-ANS encode (‘push’) and decode (‘pop’) methods.

3.2.1 The extra information

Extra information is required to initialize the bits-back chain. In practical situations, there may not be any other information that we wish to communicate, apart from that which we are modelling with a latent variable model. In this situation, we can simply send random information at the start of the chain. This means that the message has some redundancy, and the minimum amount of redundant information which we must transmit scales with the dimensionality of the latent variables. In the experiments which we present in this chapter, we use small models with a latent dimension of 40 or 50, and find that in this case the overhead is low enough that we can achieve good compression even when compressing short sequences of 30 images (see fig. 3.1). In Chapter 5, we use far larger, hierarchical models, with latent dimensionality in the hundreds of thousands, and this becomes a serious issue. In Section 5.1.2 we describe a straightforward method to deal with the overhead, by using a hybrid codec: we

first encode images using FLIF, which doesn't require extra bits, until enough of a buffer has been built up, then when possible we switch to BB-ANS, which achieves a better bit-rate.

3.2.2 Discretizing a continuous latent space

Bits back coding has previously been implemented only for models with discrete latent variables, in Frey (1997). However, many successful latent variable models utilize continuous latents, including the VAE which we use in our experiments. We present here a derivation, based on MacKay (2003), of the surprising fact that continuous latents can be coded with bits back, up to arbitrary precision, without affecting the coding rate. We also briefly discuss our implementation, which as far as we are aware is the first implementation of bits back to support continuous latents. In the following we continue to use upper case P and Q to denote mass functions for discrete distributions, and use lower case p and q for density functions of continuous distributions.

We can crudely approximate a continuous probability distribution, with density function p , with a discrete distribution, by partitioning the real line into 'buckets' of equal width δz . Indexing the buckets with $i \in I$, we assign a probability mass to each bucket of $P(i) \approx p(y_i)\delta z$, where y_i is some point in the i^{th} bucket (say its centre).

During bits back coding, we must discretize both the prior and the approximate posterior using the *same set of buckets*. Sampling from the discrete approximation $Q(i | x)$ uses approximately $\log(q(y_i | x)\delta z)$ bits, and then encoding according to the discrete approximation to the prior P costs approximately $\log(p(y_i)\delta z)$ bits. The expected message length increase for bits back with a discretized latent is therefore

$$\Delta L = -\mathbb{E}_{Q(i|x)} \left[\log \frac{P(x | y_i)p(y_i)\delta z}{q(y_i | x)\delta z} \right]. \quad (3.3)$$

The δz terms cancel, and thus the only cost to discretization results from the discrepancy between our approximation and the true, continuous, distribution. However, if the density functions are smooth (as they are in a VAE), then for small enough δz the effect of discretization will be negligible.

Note that the number of bits required to generate the latent sample scales with the precision $-\log \delta z$, meaning reasonably small precisions should be preferred in practice. Furthermore, the benefit from increasing latent precision past

a certain point is negligible for most machine learning model implementations, since they operate at 32 bit precision. In our experiments we found that increases in performance were negligible past 16 bits per latent dimension.

In our implementation, we divide the latent space into buckets which have equal mass under the prior (as opposed to equal width). This discretization is simple to implement, the computational complexity does not increase with the precision of the number of discrete intervals, and it is efficient, for the following reasons: firstly, because the discretized prior reduces to a discrete uniform distribution, which is cheap to encode/decode; secondly, because the posterior we use is parametric (a Gaussian), we can do all the computation necessary for encoding and decoding without enumerating all intervals; and thirdly (and more subtly), we don't have to worry about encoding discrete intervals which have zero mass. This is because the prior is uniform (so all intervals have the same, nonzero mass), and the only intervals we ever need to encode with the discretized posterior have been *sampled* from the discretized posterior, and thus cannot have zero mass. This saves us from a costly process of ensuring that no symbols have zero mass. Note that this discretization method can only be applied when dimensions of a latent vector are independent random variables. In Section 5.1.3 we show how to extend this simple discretization to latents with non-trivial dependence structure.

3.2.3 The need for ‘clean’ bits

In our description of bits back coding in Section 3.1, we noted that the ‘extra information’ required to seed bits back should take the form of ‘random bits’. More precisely, we need the result of mapping these bits through our decoder to produce a true sample from the distribution $q(z|x)$. A sufficient condition for this is that the bits are i.i.d. Bernoulli distributed with probability $\frac{1}{2}$ of being in each of the states 0 and 1¹. We refer to such bits as ‘clean’.

During chaining, we effectively use each compressed data point as the seed for the next. Specifically, we use the bits at the top of the ANS stack, which are the result of coding the previous latent z according to the prior $p(z)$. Will these bits be clean? The latent z is originally generated as a sample from $q(z|x)$. This distribution is clearly not equal to the prior, except in degenerate cases, so naively we wouldn't expect encoding z according to the prior to produce clean

¹This is sufficient because of a fundamental property of ANS; that integers which decode to a symbol x are distributed within the natural numbers with density $P(x)$.

bits. However, the true sampling distribution of z is in fact the *average* of $q(z|x)$ over the data distribution. That is, $q(z) \triangleq \sum_x q(z|x)P_{\text{data}}(x)$. This is referred to in Hoffman and Johnson (2016) as the ‘average encoding distribution’.

If q is equal to the true posterior, and the model is perfectly fit to the data, then $q(z) \equiv p(z)$, however in general neither of these conditions are met. Hoffman and Johnson (2016) measure the discrepancy empirically using what they call the ‘marginal KL divergence’ $D_{\text{KL}}(q(z) \parallel p(z))$, showing that this quantity contributes significantly to the ELBO for three different VAE like models learned on MNIST. This difference implies that the bits at the top the ANS stack after encoding a sample with BB-ANS will not be perfectly clean, which could adversely impact the coding rate. However, empirically we have not found this effect to be significant.

3.3 Experiments

3.3.1 Using a VAE as the latent variable model

We now demonstrate the BB-ANS coding scheme using a VAE. This model has a multidimensional latent with standard Gaussian prior and diagonal Gaussian approximate posterior:

$$p(z) = N(z; 0, I) \tag{3.4}$$

$$q(z|x) = N(z; \mu(x), \text{diag}(\sigma^2(x))). \tag{3.5}$$

We choose an output distribution (likelihood) $P(x|z)$ suited to the domain of the data we are modelling (see below). The usual VAE training objective is the ELBO, which, as we noted in Section 3.1.1, is the negative of the expected message length with bits back coding. We can therefore train a VAE as usual and plug it into the BB-ANS framework without modification.

3.3.2 Compressing MNIST

We consider the task of compressing the MNIST dataset (Lecun et al., 1998). We first train a VAE on the training set and then compress the test set using BB-ANS with the trained VAE. The MNIST dataset has pixel values in the range of integers 0, ..., 255. As well as compressing the raw MNIST data, we also present results for stochastically binarized MNIST (Salakhutdinov and Murray, 2008). For both tasks we use VAEs with fully connected generative and recognition networks, and ReLU activations.

		Binarized MNIST	Full MNIST
Raw		1	8
<i>Generic</i>	bz2	0.25	1.42
	gzip	0.33	1.64
<i>Image</i>	PNG	0.78	2.79
	WebP	0.44	2.10
<i>Models</i>	Locally masked PixelCNN ²	(0.14)	(0.65)
	Our small VAE	(0.19)	(1.39)
<i>BB-ANS</i>	Our small VAE	0.19	1.41

Table 3.3: Compression rates on the binarized MNIST and full MNIST test sets, using BB-ANS and other benchmark compression schemes, measured in bits per dimension. Note that PNG and WebP are included for context but the comparison is not particularly fair, because the image files contain metadata, the size of which is significant relative to the size of an MNIST image, particularly in the binarized case. We also give (in brackets) the negative ELBO value for our trained VAEs, and the log-likelihood under the current state of the art model, ‘Locally masked PixelCNN’, all evaluated on the test set.

For binarized MNIST the generative and recognition networks each have a single deterministic hidden layer of dimension 100, with a stochastic latent of dimension 40. The generative network outputs logits parametrizing a Bernoulli distribution on each pixel. For the full (non-binarized) MNIST dataset each network has one deterministic hidden layer of dimension 200 with a stochastic latent of dimension 50. The output distributions on pixels are modelled by a beta-binomial distribution, which is a two parameter discrete distribution. The generative network outputs the two beta-binomial parameters for each pixel.

We initialize the BB-ANS chain with a supply of ‘clean’ bits. We find that around 400 bits are required for this in our experiments. The precise number of bits required to start the chain depends on the entropy of the discretized approximate posterior (from which we are initially sampling), and scales roughly linearly with the dimensionality of the latents.

We report the achieved compression against a number of benchmarks in Table 3.3. Despite the relatively small network sizes and simple architectures we have used, the BB-ANS scheme outperforms benchmark compression schemes. While it is encouraging that even a relatively small latent variable model can outperform standard compression techniques when used with BB-ANS, the more important observation to make from Table 3.3 is that the achieved compression rate is very close to the value of the negative test ELBO seen at the end of VAE training.

In particular, the detrimental effects of finite precision, the extra information overhead (section 3.2.1), discretizing the latent (Section 3.2.2) and of less ‘clean’ bits (Section 3.2.3) do not appear to be significant. Their effects can be seen in Figure 3.4, accounting for the small discrepancy of around 1% between the negative ELBO and the achieved compression.

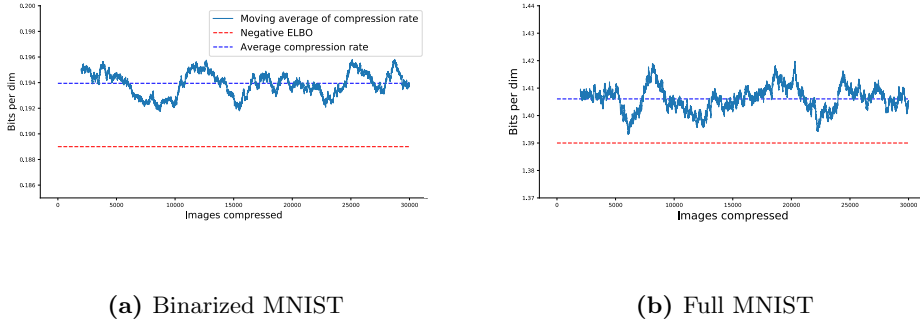


Figure 3.4: A 2000 point moving average of the compression rate, in bits per dimension, during the compression process using BB-ANS with a VAE. We compress a concatenation of three shuffled copies of the MNIST test set.

3.4 Discussion

3.4.1 Extending BB-ANS to state-of-the-art latent variable models

Implementing a state-of-the-art latent variable model is not the focus of this chapter. However, we have shown that BB-ANS can compress data to sizes very close to the negative ELBO for small-scale models. In Chapter 5 we demonstrate BB-ANS on a large-scale model and achieve state-of-the-art lossless compression on images from the ImageNet dataset.

Another extension of BB-ANS is to latent Gaussian state space models such as those studied in Johnson et al. (2016), and to state space models more generally. Very recent work (which originated from discussions during the viva for this PhD) has shown how to do this by interleaving push/pop steps with the time-steps in a model (Townsend and Murray, 2021).

3.4.2 Parallelization of BB-ANS

Modern machine learning models are optimized to exploit batch-parallelism and model-parallelism and run fastest on GPU hardware. Almost all of the computation in the BB-ANS algorithm could be executed in parallel, on GPU

²Jain et al. (2020).

hardware, and the arithmetic operations required for ANS coding should not be a computational bottleneck in a system which uses a deep generative model with ANS. In Chapter 4 we give more detail on how to implement ANS on parallel hardware, and present our own fast CPU implementation, written in NumPy (Oliphant, 2015), based on ideas from Giesen (2014).

3.4.3 Communicating the model

A neural net based model such as a VAE may have many thousands of parameters. Although it is not the focus of this work, the cost of communicating and storing a model’s parameters would need to be considered when developing a system which uses BB-ANS with a large scale model.

However, we can amortize the one-time cost of communicating the parameters over the size of the data we wish to compress. If a latent variable model could be trained such that it could model a wide class of images well, then BB-ANS could be used in conjunction with such a model to compress a large number of images. This makes the cost of communicating the model weights worthwhile to reap the subsequent gains in compression. Efforts to train latent variable models to be able to model such a wide range of images are currently of significant interest to the machine learning community, for example on expansive datasets such as ImageNet (Deng et al., 2009). We therefore anticipate that this is the most fruitful direction for practical applications of BB-ANS.

We also note that there have been many recent developments in methods to decrease the space required for neural network weights, without hampering performance. For example, methods involving quantizing the weights to low precision (Han et al., 2016; Ullrich et al., 2017), sometimes even down to single bit precision (Hubara et al., 2016), are promising avenues of research that could significantly reduce the cost of communicating and storing model weights.

3.5 Conclusion

Probabilistic modelling of data is a highly active research area within machine learning. Given the progress within this area, it is of interest to study the application of probabilistic models to lossless compression. Indeed, if practical lossless compression schemes using these models can be developed then there is the possibility of significant improvement in compression rates over existing methods.

In this chapter we have shown the existence of a scheme, BB-ANS, which can

be used for lossless compression using latent variable models. We demonstrated BB-ANS by compressing the MNIST dataset, achieving compression rates superior to generic algorithms. We have shown how to handle the issue of latent discretization. Crucially, we were able to compress to sizes very close to the negative ELBO for a large dataset. This is the first time this has been achieved with a latent variable model, and suggests that state-of-the-art latent variable models could be used in conjunction with BB-ANS to achieve significantly better lossless compression rates than current methods. All components of BB-ANS are readily parallelizable, and in Chapter 4, we describe a method for vectorizing ANS, before scaling BB-ANS up to large models and full-size colour photographs in Chapter 5.

Chapter 4

Vectorizing ANS with Craystack

To implement the prototype compression system presented in Chapter 5, we wrote new software for performing vectorized ANS operations, using NumPy (Oliphant, 2015). We call the software tool which we wrote ‘Craystack’, and make it available open source at github.com/j-towns/craystack (mirrored at doi.org/10.5281/zenodo.4707276). Tom Bird and Julius Kunze both assisted with the implementation of Craystack; Tom also created fig. 4.3. This chapter gives low-level detail on how ANS vectorization works, and a high-level overview of Craystack, as well as discussion of future directions for compression prototyping software.

We begin by describing how ANS vectorization works. The method is based on Giesen (2014). You will need to have read Chapter 2 to understand Section 4.1.

4.1 Vectorizing asymmetric numeral systems

We now describe a method for ‘vectorizing’ ANS, that is, generalizing the `push` and `pop` operations and expressing them in terms of vector or ‘single-instruction-multiple-data’ (SIMD) functions. What we give here is a high-level description, with pseudo-code. For a full implementation, which uses NumPy, see the code repository linked above.

4.1.1 The vectorized message data structure

The data structure for the message in vectorized ANS is identical to the data structure described in Section 2.2.1, except that instead of a scalar s we use a vector $s = (s_1, s_2, \dots, s_K)$. We refer to K as the ‘size’ of the message (not to be confused with the message *length*). A diagram of a vectorized message is shown in fig. 4.1. In our implementation we use a NumPy array to represent s .

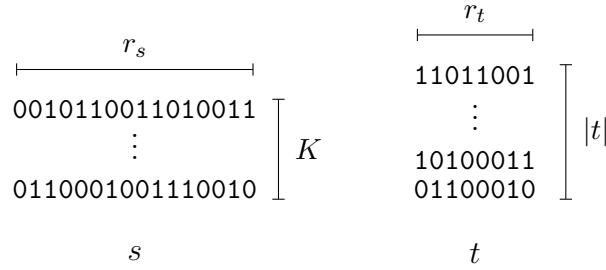


Figure 4.1: The two components of a (vectorized) message: the vector of unsigned integers s (with $r_s = 16$) and the stack of unsigned integers t (with $r_t = 8$). The integers are represented here in base 2 (binary).

4.1.2 Vectorizing the pop operation

The original (scalar) version of the `pop` function was defined in Section 2.2.2, as

```
def pop(m):
    s, t := m
    s', x := d(s)
    s, t := renorm(s', t)
    return (s, t), x
```

we needn't alter this high level definition, but just need to look at the functions `d` and `renorm`. The original definition of `d` is

```
def d(s):
    s_bar := s mod 2^r
    x, c, p := f_X(s_bar)
    s' := p * (s div 2^r) + s_bar - c
    return s', x
```

This is actually trivial to vectorize, simply by replacing all scalar arithmetic operations by their (element-wise) vector counterparts. We do not need to change the definition in any way, as long as we use arithmetic operators which are compatible with vector arguments. Note that our vectorized `d` function applied in this way will produce a vector $x = (x_1, \dots, x_K)$, and the resulting s' will also be a vector of length K .

The scalar function `renorm` was defined in Section 2.2.2 as

```
def renorm(s, t):
    while s < 2^{r_s - r_t}:
        t, t_top := stack_pop(t)
        s := 2^{r_t} * s + t_top
    return s, t
```

For vector s this can be generalized to

```
def renorm(s, t):
```

```

for  $s_k$  in  $s$ :
    while  $s_k < 2^{r_s - r_t}$ :
         $t, t_{\text{top}} := \text{stack\_pop}(t)$ 
         $s_k := 2^{r_t} \cdot s_k + t_{\text{top}}$ 
return  $s, t$ 

```

Implementing this function and its inverse, in NumPy, is difficult, so we make the simplifying assumption that $r_s \leq 2r_t$, where r_s and r_t are the precisions of the unsigned integers s and t , respectively; as in Chapter 2. This assumption has only a negligible effect on the compression rate at the precisions we typically use ($r_s = 64$ and $r_t = 32$), and implies that at most one iteration of the inner while loop will ever be required for each element, so the definition of `renorm` becomes

```

def renorm( $s, t$ ):
    for  $s_k$  in  $s$ :
        if  $s_k < 2^{r_s - r_t}$ :
             $t, t_{\text{top}} := \text{stack\_pop}(t)$ 
             $s_k := 2^{r_t} \cdot s_k + t_{\text{top}}$ 
    return  $s, t$ 

```

replacing `while` with `if`. That version of `renorm` is straightforward to implement using the NumPy `where` function and boolean indexing.

The derivation of a vectorized `push` function by inverting `pop` is mechanical, and rather than detailing it here we refer the interested reader to the implementation in the file `craystack/rans.py` in the Craystack repository.

4.1.3 The length of the vectorized message

In Section 2.2.1, we defined the *length* of a scalar message as

$$l(m) := r_s + r_t |t| \tag{4.1}$$

and the *effective length* as

$$l^*(m) := \log s + r_t |t|. \tag{4.2}$$

The length of the scalar message tells us how many bits are required to store a flattened representation of that message. We could flatten a vector message $m = ((s_1, \dots, s_K), t)$, simply by concatenating the elements s_1, \dots, s_K and the

elements of t , which would lead to a flattened string with length

$$l_{\text{naïve}}(m) = Kr_s + r_t|t|. \quad (4.3)$$

We can improve on this length, by approaching the question of how to flatten m as a *lossless compression problem*. We split the vector s into a pair $s_1, (s_2, \dots, s_K)$, and treat $m' = (s_1, t)$ as a *scalar* message, and the question becomes, how best to *push* the elements s_2, \dots, s_K onto m' ? We can use ANS (a near-optimal compressor) to do this, and then we just have to pick a distribution to use to model s_2, \dots, s_K .

4.1.3.1 The Benford distribution

It has been observed empirically, by Bloom (2014), that in practical settings, the elements s_k precisely follow *Benford's law* (Benford, 1938). Benford's law says that samples of leading digits from sets of 'real-world' numbers tends to follow a distribution with mass function

$$P_{\text{Benford}}(x) \propto 1/x, \quad (4.4)$$

where x is the integer formed by concatenating the r leading digits of a sampled constant, for some fixed value of r (the normalizing constant of P is a function of r). We refer to this as the 'Benford distribution'. MacKay (2003) argues that the reason that digits in physical constants must follow this distribution is that their distribution ought not to depend on the units of measurement used, and that the Benford distribution is the unique distribution which is invariant to rescaling of units.

The number s can be thought of as the first $r_s/(r_s - r_t)$ digits of a flat message, if the message is treated as a large integer, expressed in base $r_s - r_t$. Since, as we have shown, message length is approximately equal to the information content of an encoded sequence plus a constant, the fact that s is Benford distributed is equivalent to saying that *probability masses of sample sequences tend to follow Benford's law*.

Why might this be? Similarly to MacKay (2003), we can sketch a proof by contradiction. Consider the sequence

$$p_n = \prod_{i=1}^n P(x_i | x_1, \dots, x_{i-1}) \quad (4.5)$$

for a sample $x_1, x_2, \dots, x_n, \dots$ with mass function P . Let q_n be the integer composed of the first r non-zero digits of p_n . If we assume that q_n tends to a stationary distribution as n increases, and that the stationary distribution is positive (that is, each possible state of q_n has non-zero probability), then by definition the stationary distribution must be invariant to the scale of p_n , because for each n , q_{n+1} is derived from a rescaling of p_n (specifically, multiplying p_n by $P(x_{n+1} | x_1, \dots, x_n)$).

We can approximate the normalizer for the Benford distribution over s

$$\sum_{s=2^{r_s-r_t}}^{2^{r_s}} \frac{1}{s} \approx \int_{2^{r_s-r_t}}^{2^{r_s}} \frac{1}{s} \quad (4.6)$$

$$= [r_s - (r_s - r_t)] \ln 2 \quad (4.7)$$

$$= r_t \ln 2 \quad (4.8)$$

and thus, using the Benford distribution to append the elements s_2, \dots, s_K to the scalar message (s_1, t) results in a flattened message whose length satisfies

$$l_{\text{opt}}(m) \leq r_s + \sum_{k=2}^K \log s_k + (K-1) \log(r_t \ln 2) + r_t |t| + (K-1)\epsilon, \quad (4.9)$$

by applying eq. (2.16) to each s_2, \dots, s_K .

4.1.4 The cost of vectorization

In the scalar case we were able to prove the following bound on message length, eq. (2.21),

$$l(m) \leq h(x_1, \dots, x_N) + N\epsilon + r_s. \quad (4.10)$$

If we initialize a vector message to $m = (s_{\text{init}}, t_{\text{init}})$ where $s_{\text{init}} := (2^{r_s-r_t}, 2^{r_s-r_t}, \dots, 2^{r_s-r_t})$, i.e. the minimum permissible value for each element, then apply eq. (4.9) and eq. (2.16) to each encoding iteration element-wise, we get an inequality which is analogous to eq. (4.10)

$$l_{\text{opt}}(m) \leq h(x_1, \dots, x_N) + (K-1) \log(r_t \ln 2) + K(N\epsilon + r_s - r_t) + r_t. \quad (4.11)$$

For short messages and large K , the overhead terms $(K-1) \log(r_t \ln 2)$ and $K(r_s - r_t)$ may be significant.

One approach to reducing these overheads, which we use in Chapter 5, is to encode the first samples in a sequence using *scalar* ANS, then expand the size of

the message progressively by *decoding* (i.e. sampling) new elements s_k from the message as it grows in length. The $K(r_s - r_t)$ overhead term from the vector initial message disappears when we do this, and if we use the Benford distribution to sample the s_k then the $(K - 1) \log(r_t \ln 2)$ term is cancelled out as well, meaning we can vectorize with negligible overhead in compression rate. However, this approach is complicated to implement, and it is unclear whether it would be sensible to use it in practical settings given the additional complexity required.

4.1.5 The benefit of vectorization

The performance benefits of vectorization are drastic. To demonstrate this we benchmarked the ANS encoding and decoding of ImageNet images (see Chapter 5 for more details on the codec), using a scalar ANS implementation vs a vectorized implementation. The results are shown in fig. 4.2, and show that in this setting the vectorized ANS implementation was nearly three orders of magnitude faster than the scalar implementation. This is because of loops which occur in the Python interpreter in the scalar version being ‘pushed down’ into highly optimized NumPy kernels, implemented in C and Fortran.

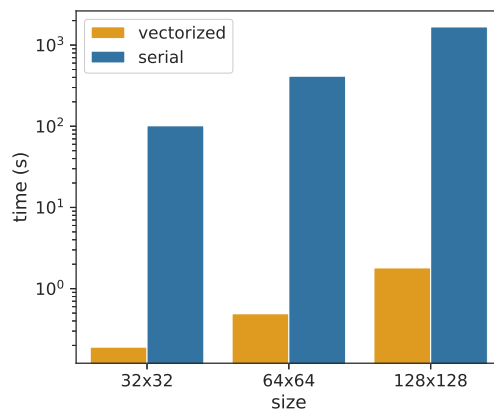


Figure 4.2: Runtime of vectorized vs. serial ANS implementations for different image sizes. Times were computed on a desktop with 6 CPU cores.

4.2 Craystack

Having described some of the low-level details of vectorized ANS, we now discuss some of the high-level features of the Craystack library, which aims to make prototyping lossless compression systems with ANS straightforward and approachable for machine learning practitioners.

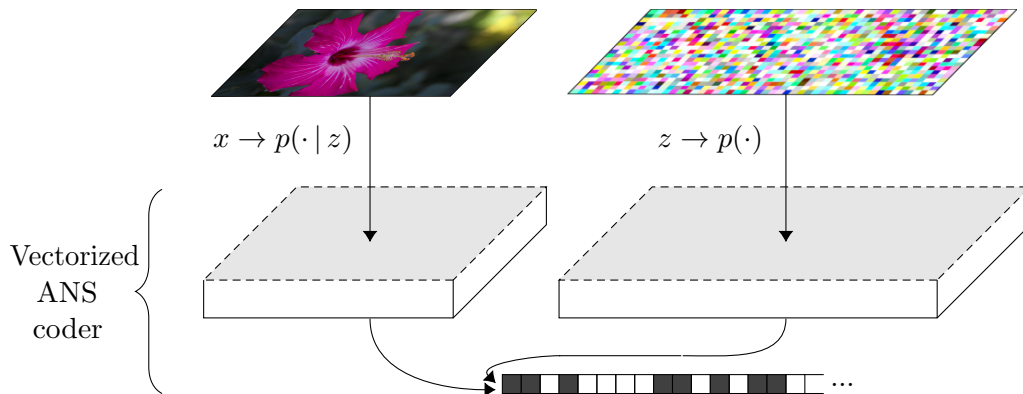


Figure 4.3: Visualizing the process of pushing images and latents from a VAE to the vectorized ANS stack with Craystack. The ANS stack head is shaped such that images and latents can be pushed and popped in parallel, without reshaping. Beneath the shaped top of the stack is the flat message stream output by ANS.

4.2.1 High-level design

When writing a lossless compression system, it is essential to ensure that the decoder function is the precise inverse of the encoder. Manually writing code which satisfies this constraint requires care, and one of the central goals of Craystack is to ease this burden by ensuring that this constraint is satisfied *automatically*. To achieve this, Craystack provides a set of primitive codecs, including codecs for coding data according to uniform, Bernoulli, discretized logistic, discretized Gaussian, and categorical distributions. Each primitive codec is comprised of a (`push`, `pop`) inverse pair. As well as the primitives, we also provide a set of combinators for composing and nesting codecs. Examples of high level combinators included in Craystack are a BB-ANS combinator and a combinator for compression with an auto-regressive model.

The code for the BB-ANS combinator is shown in fig. 4.4. By using the Craystack primitives combined with the provided combinators, users are able to build elaborate compressors without needing to worry about the correctness of their code.

4.2.2 Generalizing vectorized ANS

Another way in which Craystack aims for convenience is by allowing arbitrarily shaped data arrays to be coded directly. The ‘vector’ s in the Craystack ANS implementation is a 1D NumPy array, but we allow codecs to use arbitrary shaped NumPy views of s , and multiple views of s can even be separated into an arbitrary nesting of Python tuples, lists and dictionaries. Figure 4.3 shows a visualization of this, where s is split into a pair of rectangular views.


```

def BBANS(prior, likelihood, posterior):
    def push(message, data):
        message, latent = posterior(data).pop(message)
        message = likelihood(latent).push(message, data)
        message = prior.push(message, latent)
        return message

    def pop(message):
        message, latent = prior.pop(message)
        message, data = likelihood(latent).pop(message)
        message = posterior(data).push(message, latent)
        return message, data
    return Codec(push, pop)

```

Figure 4.4: The BB-ANS combinator provided in Craystack. It accepts codecs for the prior, likelihood and approximate posterior as arguments and returns a codec which uses BB-ANS to compress data modelled with a latent variable model.

The API for this mechanism, which is via a function called `substack`, only requires the user to supply a function which creates the view from a flat vector s . For example, codecs for the latent and observation in fig. 4.3 could be expressed as

```

latent_view = lambda s: s[:latent_size].reshape(latent_shape)
obs_view    = lambda s: s[latent_size:].reshape(obs_shape)

flat_latent_codec = substack(shaped_latent_codec, latent_view)
flat_obs_codec    = substack(shaped_obs_codec,    obs_view)

```

The view functions are linear, orthonormal mappings (in the sense of linear algebra), and therefore they can be *inverted* using reverse mode automatic differentiation (AD). This is because the derivative of a linear function is the function itself; reverse mode AD computes the transpose of the derivative, and the transpose of an orthonormal mapping is equal to the mapping's (pseudo) inverse. The `substack` function automatically computes the inverse of the view functions using Autograd (Maclaurin et al., 2015), ensuring that the `push` and `pop` functions produced are inverse to one-another. Methods for implementing automatic view inversion from scratch (i.e. without exploiting an AD system), can be found in Voigtländer (2009).

4.3 Future directions for Craystack

We broadly see two major directions for improving Craystack. The first is to allow ANS to run on accelerators such as GPUs. One way to achieve this

would be to port the implementation into Google’s JAX software (Bradbury et al., 2018), which enables use of GPU and TPU backends with a NumPy-like API. JAX also allows just-in-time compilation with automatic loop fusion, amongst other optimizations, and implementing deep generative models in JAX is reasonably straightforward¹.

One interesting open question is whether it is possible to have a user write a decoder (`pop`) function (which would be very similar to a function for sampling from a model), and to automatically transform that function into an encoder (`push`) function. Internally, JAX is based on a core system for writing composable function transformations, with transformations such as forward and reverse mode differentiation, auto-vectorization and just-in-time compilation provided. More investigation is needed to find out whether this transformation can be implemented in a practical manner, but if so it could fit well into the existing JAX system.

¹During the very recent work described in Ruan et al. (2021), we implemented a proof-of-concept version of Craystack’s core vectorized rANS in JAX, and have run it on GPU, though we haven’t thoroughly benchmarked it yet. The implementation is available at github.com/j-towns/crayjax, mirrored at <https://doi.org/10.5281/zenodo.4650348>.

Chapter 5

Scaling up bits back coding with asymmetric numeral systems

This chapter is based on the paper ‘HiLLoC: lossless image compression with hierarchical latent variable models’, published at ICLR 2020 (Townsend et al., 2020). We show that the bits back with ANS method presented in Chapter 3 can be scaled up to larger models and applied to compression of colour photographs, achieving a state-of-the-art compression rate on full-sized images from the ImageNet dataset (Russakovsky et al., 2015). This was a collaboration with Tom Bird, who contributed to the ideas and wrote the paper and experiments with me, and Julius Kunze, who mostly assisted in writing the experiments.

In order to scale up the methods in Chapter 3, we will use four novel ideas:

1. A vectorized ANS implementation supporting dynamic shape.
2. Direct coding of arbitrary sized images using a fully convolutional model.
3. Dynamic discretization of hierarchical latents.
4. Initializing the bits back chain using a different codec.

We have already discussed item 1 in detail in Chapter 4, and will discuss items 2 to 4 in Section 5.1. We call the combination of BB-ANS using a hierarchical latent variable model and the above techniques: ‘Hierarchical Latent Lossless Compression’ (HiLLoC). In our experiments (Section 5.2), we demonstrate that HiLLoC can be used to compress colour images from the ImageNet test set at rates close to the ELBO, outperforming all of the other codecs which we benchmark.

Note that the ‘Bit-Swap’ method, presented by F. H. Kingma et al. (2019) has a similar aim to HiLLoC, namely scaling up the method described in

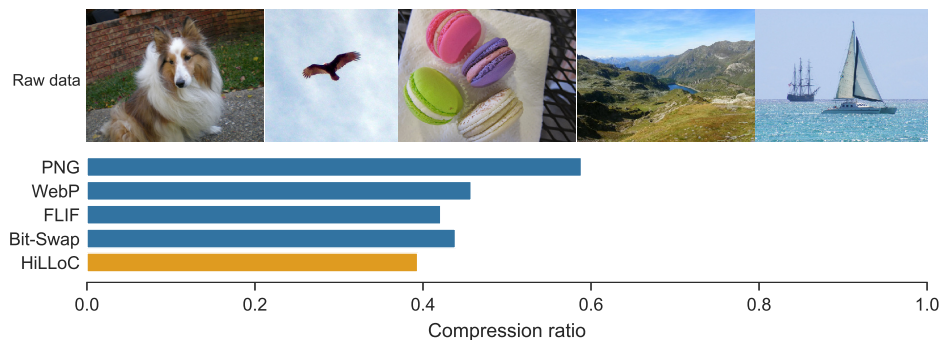


Figure 5.1: A selection of images from the ImageNet dataset and the compression rates achieved on the dataset by PNG, WebP, FLIF, Bit-Swap and the HiLLoC codec (with ResNet VAE) presented in this chapter.

Chapter 3. We describe the central idea of Bit-Swap in Section 5.1.2.1, and where relevant we compare our approach to it, describing the trade-offs where they exist. We recommend reading at least Chapters 2 and 3 for background before trying to understand the material in this chapter.

5.1 Scaling up BB-ANS

5.1.1 Fully convolutional models

When all of the layers in the generative and recognition networks of a VAE are either convolutional or element-wise functions (i.e. the VAE has no densely connected layers), then it is possible to evaluate the recognition network on images of any height and width, and similarly to pass latents of any height and width through the generative network to generate an image. Thus, such a VAE can be used as a (probabilistic) model for images of any size.

We exploit this fact, and show empirically in Section 5.2 that, surprisingly, a fully convolutional VAE trained on 32×32 images can perform well (in the sense of having a high ELBO) as a model for 64×64 images as well as far larger images. This, in turn, corresponds to a good compression rate, and we implement lossless compression of arbitrary sized images by using a VAE in this way. We would expect this result might extend to deep generative models other than VAEs, i.e. models which use masked convolutions (Germain et al., 2015) and flow-based models, although investigating this is outside of the scope of this work.

5.1.2 Reducing one-time overhead

As mentioned in Section 3.1.2, if we don't have an 'extra message' to communicate, and we wish to communicate a sequence of data points (modeled

as independent), then we may use a method other than BB-ANS to code the first data, until enough of a buffer has been built up to generate a latent from $Q(z|x)$, or we may generate a latent z for the first datapoint, using any process we like, and encode it with the prior and posterior, accepting that we will then be communicating redundant information at the beginning of a chain. The overhead of starting the chain is particularly significant for hierarchical models, because the information required to sample from $Q(z|x)$ typically scales with the dimensionality of z and thus with the latent depth.

In this subsection we discuss two techniques for mitigating this issue. The first is ‘Bit-Swap’, which was presented in F. H. Kingma et al. (2019). Bit-Swap exploits Markov chain structured latents, by interleaving latent pop and push steps, and for models with the right structure it can reduce the overhead from $O(L)$ to $O(1)$.

It may sometimes be desirable to use latents which do not have a Markov chain structure, and for this we use an alternative, simpler method for reducing overhead (which could also be used to *further* reduce the overhead when using Bit-Swap), which is simply to use the FLIF codec, which we know to be reasonably efficient, to code the first elements of the batch.

5.1.2.1 Bit-Swap

The Bit-Swap method (F. H. Kingma et al., 2019) assumes that the generative model has a hierarchical structure with L latent layers, illustrated, for a model with $L = 3$, in fig. 5.2. In particular it assumes that the generative joint distribution and the approximate posterior may be expressed as

$$P(x, z) = P(x | z_1) \left[\prod_{l=1}^{L-1} P(z_l | z_{l+1}) \right] P(z_L) \quad (5.1)$$

$$Q(z | x) = Q(z_1 | x) \prod_{l=2}^L Q(z_l | z_{l-1}). \quad (5.2)$$

Given this structure, it is possible to alter the BB-ANS procedure from table 3.2, to avoid popping the whole of z in one go. The sender can interleave steps, popping one layer at a time and pushing whenever possible, as shown in table 5.1. This technique reduces the initial bits overhead from $O(L)$ to $O(1)$, at the cost of introducing the Markov restriction on the generative model and approximate posterior (we discuss this further in Section 5.1.3.1).

Variables	Operation	Operation	Variables
x	$z_1 \leftarrow Q(\cdot x)$	$z_L \leftarrow P(\cdot)$	z_L
x, z_1	$x \rightarrow P(\cdot z_1)$	$z_{L-1} \leftarrow P(\cdot z_L)$	z_{L-1}, z_L
z_1	$z_2 \leftarrow Q(\cdot z_1)$	\vdots	\vdots
z_1, z_2	$z_1 \rightarrow P(\cdot z_2)$	$z_1 \leftarrow P(\cdot z_2)$	z_1, z_2
\vdots	\vdots	$z_2 \rightarrow Q(\cdot z_1)$	z_1
z_{L-1}, z_L	$z_{L-1} \rightarrow P(\cdot z_L)$	$x \leftarrow P(\cdot z_1)$	x, z_1
z_L	$z_L \rightarrow P(\cdot)$	$z_1 \rightarrow Q(\cdot x)$	x

(a) Encoder
(b) Decoder

Table 5.1: The Bit-Swap encoder and decoder procedures, in order from the top, for a model with L layers. For the encoder, the ‘Variables’ column shows variables known before each operation. For the decoder, that column shows variables known after each operation.

5.1.2.2 Reducing overheads with FLIF

As will be explained in Section 5.1.3, our dynamic discretization method precludes the use of Bit-Swap for reducing the one-time cost of starting a BB-ANS chain, and we also want to be able to use a model without Markov chain latent structure. We propose instead to use a significantly simpler method to address the high cost of coding a small number of samples with BB-ANS: we code the first samples using a different codec. The purpose of this is to build up a sufficiently large buffer of compressed data to permit the first stage of the BB-ANS algorithm—to pop a latent sample from the posterior. In our experiments we use the ‘Free Lossless Image Format’ (FLIF; Sneyers and Wuille, 2016) to build up the buffer. We chose this codec because it was the best performing at the time of writing, but in principal any lossless codec could be used.

The amount of previously compressed data required to pop a posterior sample from the ANS stack (and therefore start the BB-ANS chain) is roughly proportional to the size of the image we wish to compress, since in a fully convolutional model the dimensionality of the latent space is determined by the image size.

We can exploit this to obtain a better compression rate than FLIF as quickly as possible. We do so by partitioning the first images we wish to compress with HiLLoC into smaller patches. These patches require a smaller data buffer, and thus we can use the superior HiLLoC coding sooner than if we attempted to compress full images. We find experimentally that, generally, larger patches have

a better coding rate than smaller patches. Therefore we increase the size of the image patches being compressed with HiLLoC as more images are compressed and the size of the data buffer grows, until we finally compress full images directly once the buffer is sufficiently large.

For our experiments on compressing full ImageNet images, we compress 32×32 patches, then 64×64 , then 128×128 before switching to coding the full size images directly. Note that since our model can compress images of any shape, we can compress the edge patches which will have different shape if the patch size does not divide the image dimensions exactly. Using this technique means that our coding rate improves gradually from the FLIF coding rate towards the coding rate achieved by HiLLoC on full images. We compress only 5 full ImageNet images using FLIF before we are able to start compressing 32×32 patches using HiLLoC.

5.1.3 Dynamic discretization

It is standard for state of the art latent variable models to use continuous latent variables. Since ANS operates over *discrete* probability distributions, if we wish to use BB-ANS with such models it is necessary to discretize the latent space so that latent samples can be communicated. In Section 3.2.2, we described a *static* discretization scheme for the latents in a simple VAE with a single layer of continuous latent variables, and in Section 3.3, we showed that this discretization has a negligible impact on compression rate. The addition of multiple layers of stochastic variables to a VAE has been shown to improve performance (D. P. Kingma et al., 2016; Sønderby et al., 2016; Maaløe et al., 2019; F. H. Kingma et al., 2019). Motivated by this, we propose a discretization scheme for hierarchical VAEs with multiple layers of latent variables.

The discretization described in Section 3.2.2 is formed by dividing the latent space into intervals of equal mass under the prior $p(z)$. For a hierarchical model, the prior on each layer depends on the previous layers:

$$p(z_{1:L}) = p(z_L) \prod_{l=1}^{L-1} p(z_l | z_{l+1:L}). \quad (5.3)$$

It isn't immediately possible to use the simple static scheme from Section 3.2.2, since the marginals $p(z_1), \dots, p(z_{L-1})$ are not known. The Bit-Swap method, described in Section 5.1.2.1, estimates these marginals by sampling, creating static bins based on the estimates. They demonstrate that this approach can

work well, see F. H. Kingma et al. (2019) for details. We propose an alternative approach, allowing the discretization to vary with the context of the latents we are trying to code. We refer to this as *dynamic discretization*.

In dynamic discretization, instead of discretizing with respect to the marginals of the prior, we discretize according to the *conditionals* in the prior, $p(z_l | z_{l+1:L})$. Specifically, for each latent layer l , we partition each dimension into intervals which have equal probability mass under the conditional $p(z_l | z_{l+1:L})$. This directly generalizes the static scheme from Section 3.2.2.

Dynamic discretization is more straightforward to implement than the method used with Bit-Swap, because it doesn't require calibrating the discretization to samples. However it imposes a restriction on model structure, in particular it requires that posterior inference is done *top-down*. This precludes the use of the Bit-Swap technique for reducing the size of the initial message needed to start the bits-back chain. In Section 5.1.3.1 we contrast the model restriction from dynamic discretization with the bottom-up, Markov restriction imposed by Bit-Swap itself.

We give further details about the dynamic discretization implementation which we use in Appendix B.

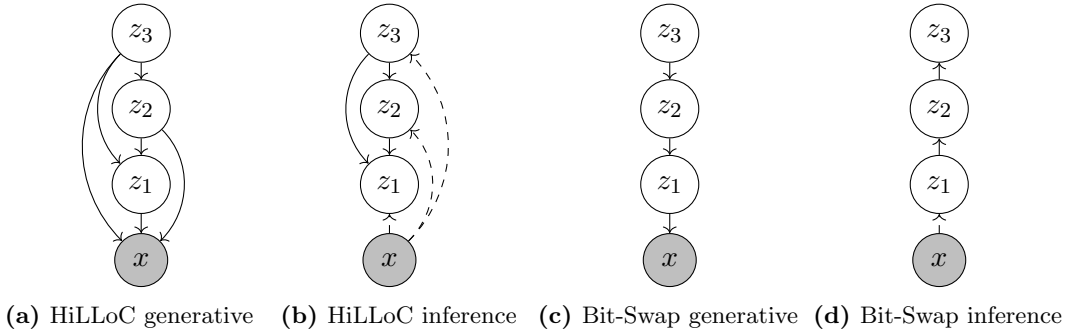


Figure 5.2: Graphical models representing the generative and inference models with HiLLoC and Bit-Swap, both using a 3 layer latent hierarchy. The dashed lines indicate dependence on a fixed observation.

5.1.3.1 Model restrictions

The first stage of BB-ANS encoding is to pop from the posterior, $z_{1:L} \leftarrow q(\cdot | x)$. When using dynamic discretization, popping the layer z_l requires knowledge of the discretization used for z_l and thus of the conditional distribution $p(z_l | z_{l+1:L})$. This requires the latents $z_{l+1:L}$ to have already been popped. Because of this, latents in general must be popped (sampled) in ‘top-down’ order, i.e. z_L first, then z_{L-1} and so on down to z_1 .

The most general form of posterior for which top-down sampling is tractable is

$$q(z_{1:L} | x) = q(z_L | x) \prod_{l=1}^{L-1} q(z_l | z_{l+1:L}, x). \quad (5.4)$$

This is illustrated, for a hierarchy of depth 3, in fig. 5.2b. The Bit-Swap technique (F. H. Kingma et al., 2019) requires that inference be done bottom up, and that generative and inference models must both be a Markov chain on z_1, \dots, z_L , and thus cannot use skip connections. These constraints are illustrated in figs. 5.2c and 5.2d. Skip connections have been shown to improve the achievable ELBO value in very deep models (Sønderby et al., 2016; Maaløe et al., 2019). HiLLoC does not have this constraint, and we do utilize skip connections in our experiments.

5.2 Experimental results

Using Craystack, we implement HiLLoC with a ResNet VAE (RVAE; D. P. Kingma et al., 2016). In all experiments we used an RVAE with 24 stochastic hidden layers. The RVAE utilizes skip connections, which are important for effectively training models with such a deep latent hierarchy. See Appendix C for more details.

We trained the RVAE on the ImageNet 32 training set, then evaluated the RVAE ELBO and HiLLoC compression rate on the ImageNet 32 test set. To test generalization, we also evaluated the ELBO and compression rate on the test sets of ImageNet64, CIFAR-10 and full size ImageNet. For full size ImageNet, we used the partitioning method described in Section 5.1.2 to compress the first images. The results are shown in Table 5.2.

For HiLLoC the compression rates are for the entire test set, except for full ImageNet, where we use 2000 randomly selected images from the test set.

Table 5.2 shows that HiLLoC achieves competitive compression rates on all benchmarks, and is the first deep learning based method to be evaluated on full size ImageNet images. This significant upscaling relative to previous work was enabled by the speedups from vectorizing ANS in Craystack. We anticipate that flow-based and/or autoregressive models may outperform VAEs on this task

¹Integer discrete flows, retrieved from Hoogeboom et al. (2019).

²Integer discrete flows trained on ImageNet 32. ImageNet 64 images are split into four 32×32 patches. Retrieved from Hoogeboom et al. (2019).

³Local bits back, retrieved from Ho et al. (2019).

⁴For Bit-Swap, full size ImageNet images were cropped so that their side lengths were multiples of 32.

Table 5.2: Compression performance of HiLLoC with RVAE compared to other codecs. Rates measured in bits/dimension (raw data is 8 bits/dimension). For HiLLoC we display compression rate and theoretical performance (ELBO). All HiLLoC results are obtained from the *same model*, trained on ImageNet 32.

		ImageNet 32	ImageNet 64	Cifar-10	ImageNet
<i>Generic</i>	PNG	6.39	5.71	5.87	4.71
	WebP	5.29	4.64	4.61	3.66
	FLIF	4.52	4.19	4.19	3.37
<i>Flow-based</i>	IDF ¹	4.18	3.90	3.34	-
	IDF generalized ²	4.18	3.94	3.60	-
	LBB ³	3.88	3.70	3.12	-
<i>VAE-based</i>	Bit-Swap	4.50	-	3.82	3.51 ⁴
	HiLLoC	4.20	3.90	3.56	3.15
	HiLLoC (ELBO)	(4.18)	(3.89)	(3.55)	(3.14)

in terms of compression rate, possibly with slower compression/decompression times. Further work is required to investigate different model classes and architectures.

The fact that HiLLoC with an RVAE can achieve state of the art compression on full ImageNet relative to the baselines, even under a change of distribution, is striking. This provides evidence of its efficacy as a general method for lossless compression of natural images. Naïvely, one might expect a degradation of performance relative to the original test set when changing the test distribution—even more so when the resolution changes. However, in the settings we studied, the *opposite* was true, in that the average per-pixel negative ELBO (and thus the compressed message length) was *lower* on all other datasets compared to the ImageNet 32 validation set.

In the case of CIFAR, we conjecture that the reason for this is that its images are simpler and contain more redundancy than ImageNet. This theory is backed up by the performance of standard compression algorithms which, as shown in Table 5.2, also perform better on CIFAR images than they do on ImageNet 32. We find the compression rate improvement on larger images more surprising. We hypothesize that this is because pixels at the edge of an image are harder to model because they have less context to reduce uncertainty. The ratio of edge pixels to interior pixels is lower for larger images, thus we might expect less uncertainty per pixel in a larger image.

To demonstrate the effect of vectorization we timed ANS of single images at different, fixed, sizes, using a fully vectorized and a fully serial implementation. The results are shown in Figure 4.2 in the previous chapter, which clearly shows

a speed-up of nearly three orders of magnitude for all image sizes. We find that the run times for encoding and decoding are roughly linear in the number of pixels, and the time to compress/decompress an average sized ImageNet image of 500×374 pixels (with vectorized ANS) is around 30s on a desktop computer with 6 CPU cores and a GTX 1060 GPU. Most of this time is spent on neural net computation, highlighting the need for more computationally efficient models, and/or more powerful hardware.

5.3 Discussion

Our experiments demonstrate HiLLoC as a bridge between large scale latent variable models and compression. To do this we use simple variants of pre-existing VAE models. Having shown that bits back coding is flexible enough to compress well with large, complex models, we see plenty of work still to be done in searching model structures (i.e. architecture search), optimizing with a trade-off between compression rate, encode/decode time and memory usage. Particularly pertinent for HiLLoC is latent dimensionality, since compute time and memory usage both scale with this. Since the model must be stored/transmitted to use HiLLoC, weight compression is also highly relevant. This is a well-established research area in machine learning (Han et al., 2016; Ullrich et al., 2017).

Our experiments also demonstrated that one can achieve good performance on a dataset of large images by training on smaller images. This result is promising, but future work should be done to discover what the best training datasets are for coding generic images. One question in particular is whether results could be improved by *training*, as opposed to just evaluating, models on larger images and/or images of varying shape and size. We leave this to future work. Another related direction for improvement is batch compression of images of different sizes using masking, analogous to how samples of different length may be processed in batches by recurrent neural nets in natural language processing applications.

Whilst this work has focused on latent variable models, there is also promise in applying state of the art fully observed auto-regressive models to lossless compression. We look forward to future work investigating the performance of models such as WaveNet (van den Oord et al., 2016a) for lossless audio compression as well as PixelCNN++ (Salimans et al., 2016) and the state of the art models in Jun et al. (2020) for images. Sampling speed for these models, and

thus decompression, scales with autoregressive sequence length, and can be very slow. This could be a serious limitation, particularly in common applications where encoding is performed once but decoding is performed many times. This effect can be mitigated by using dynamic programming (Le Paine et al., 2016; Ramachandran et al., 2017), and altering model architecture (Reed et al., 2017), but on parallel architectures sampling/decompression may still be slower than with VAE models.

On the other hand, fully observed models, as well as the flow based models of Hoogeboom et al. (2019) and Ho et al. (2019), do not require bits back coding, and therefore do not have to pay the one-off cost of starting a chain. Therefore they may be well suited to situations where one or a few i.i.d. samples are to be communicated. Similar to the way that we use FLIF to code the first images for our experiments, one could initially code images using a fully observed model then switch to a faster latent variable model once a stack of bits has been built up.

5.4 Conclusion

In this chapter, we described HiLLoC, an extension of BB-ANS to hierarchical latent variable models, and show that HiLLoC can perform well with large models. We open-sourced our implementation, along with the Craystack package for prototyping lossless compression. We have also explored generalization of large VAE models, and established that fully convolutional VAEs can generalize well to other datasets, including images of very different size to those they were trained on. Finally, we have described a method for compressing images of arbitrary size with HiLLoC, achieving a compression rate superior to the best available codecs on ImageNet images.

Chapter 6

General Conclusions

The research on which this thesis is based had two aims. The first, more direct, aim was to develop methods for lossless compression using latent variable models, and show that they could perform well. The second was broadly to demonstrate that lossless compression is an interesting and potentially useful application of recently developed generative modelling techniques, and that compression ideas can be prototyped by machine learning practitioners without too much difficulty, particularly using ANS.

Towards the first aim, we have demonstrated that it is practically possible to implement lossless compression using large scale latent variable models, and in particular that compression rates very close to the model negative ELBO can be achieved in large-scale VAEs, on full size images from the ImageNet dataset. We have also, towards the second aim, presented what we hope is an accessible introduction to ANS, and easy-to-use software for prototyping ANS-based compression.

It is too early to say how significant our impact has been in shifting attention towards lossless compression among machine learning practitioners. However, there has already a strong response from the community: whilst we're not aware of any work prior to ours connecting deep generative models to ANS, since the publication of Townsend et al. (2019) there have been at least four publications by other machine learning researchers which use ANS and deep generative models to demonstrate new lossless compression ideas (F. H. Kingma et al., 2019; Ho et al., 2019; Hoogeboom et al., 2019; van den Berg et al., 2020). Perhaps more significantly, as a result of our work, ANS and BB-ANS now form part of the popular 'Deep Unsupervised Learning' course at UC Berkeley. We hope that the ideas, experimental results and software presented in this thesis can form a strong basis for future research applying deep learning to lossless compression.

Appendix A

Example ANS implementation

The following is a complete implementation, in pure Python, of the range variant of asymmetric numeral systems (rANS), using the same variable names as in Chapter 2. A cons list (implemented as a binary tuple tree) is used as the stack data structure for t , and the methods `flatten_stack` and `unflatten_stack` convert between the stack representation and a Python list. For more detail, see github.com/j-towns/ans-notes/blob/master/rans.py.

Listing A.1: Scalar rANS core

```
s_prec = 64
t_prec = 32
t_mask = (1 << t_prec) - 1
s_min = 1 << s_prec - t_prec

m_init = s_min, () # Shortest possible message

def rans(model):
    f, g, p_prec = model
    def push(m, x):
        s, t = m
        c, p = g(x)
        while s >= p << s_prec - p_prec:
            s, t = s >> t_prec, (t, s & t_mask)
        s = (s // p << p_prec) + s % p + c
        return s, t

    def pop(m):
        s, t = m
        s_bar = s & ((1 << p_prec) - 1)
        x, (c, p) = f(s_bar)
        s = p * (s >> p_prec) + s_bar - c
        while s < s_min:
            t, t_top = t
            s = (s << t_prec) + t_top
        return (s, t), x
    return push, pop
```

```

def flatten_stack(t):
    flat = []
    while t:
        t_top, t = t
        flat.append(t_top)
    return flat

def unflatten_stack(flat):
    t = ()
    for t_top in reversed(flat):
        t = t_top, t
    return t

```

Listing A.2: Scalar rANS example usage

```

import math

log = math.log2

# We encode some data using the example model from
# Chapter 2 and verify the inequality in eq. (2.19).

# First setup the model
p_prec = 3

# Probability weights, must sum to 2 ** p_prec
ps = {'a': 1,
      'b': 2,
      'c': 3,
      'd': 2}

# Cumulative probabilities
cs = {'a': 0,
      'b': 1,
      'c': 3,
      'd': 6}

# Backwards mapping
s_bar_to_x = {0: 'a',
              1: 'b', 2: 'b',
              3: 'c', 4: 'c', 5: 'c',
              6: 'd', 7: 'd'}

def f(s_bar):
    x = s_bar_to_x[s_bar]
    c, p = cs[x], ps[x]
    return x, (c, p)

def g(x):
    return cs[x], ps[x]

model = f, g, p_prec

```

```

push, pop = rans(model)

# Some data to compress
xs = ['a', 'b', 'b', 'c', 'b', 'c', 'd', 'c', 'c']

# Compute h(xs):
h = sum(map(lambda x: log(2 ** p_prec / ps[x]), xs))
print('Information content of sequence: '
      'h(xs) = {:.2f} bits.'.format(h))
print()

# Initialize the message
m = m_init

# Encode the data
for x in xs:
    m = push(m, x)

# Verify the inequality in eq (20)
eps = log(1 / (1 - 2 ** -(s_prec - p_prec - t_prec)))
print('eps = {:.2e}'.format(eps))
print()

s, t = m
lhs = (log(s) + t_prec * len(flatten_stack(t))
       - log(s_min))
rhs = h + len(xs) * eps
print('Eq (20) inequality, rhs - lhs == {:.2e}'
      .format(rhs - lhs))
print()

# Decode the message, check that the decoded data
# matches original
xs_decoded = []
for _ in range(len(xs)):
    m, x = pop(m)
    xs_decoded.append(x)

xs_decoded = reversed(xs_decoded)

for x_orig, x_new in zip(xs, xs_decoded):
    assert x_orig == x_new

# Check that the message has been returned to its
# original state
assert m == m_init
print('Decode successful!')

```


Appendix B

Reparameterizing discretized latents in hierarchical VAEs

After discretizing the latent space, the latent variable at layer l can be treated as simply an index i_l labeling the interval into which z_l falls. We introduce the following notation for pushing and popping according to a discretized version of the posterior:

$$i_l \leftrightarrow Q_l(\cdot | i_{l+1:L}, x), \quad (\text{B.1})$$

where $Q_l(\cdot | i_{l+1:L}, x)$ is the distribution over the intervals of the discretized latent space for z_l , with interval masses equal to their probability under $q(z_l | \tilde{z}_{l+1:L}, x)$. The discretization is created by splitting the latent space into equal mass intervals under $p(z_l | \tilde{z}_{l+1:L})$. We use \tilde{z} to denote the centred values that can be reconstructed from the indices i_l . To be precise, $\tilde{z}_l(i_l)$ is set to the median value within the interval indexed by i_l , under the prior. Note that Q_l has an implicit dependence on the previous prior distributions $p(z_k | z_{k+1:L})$ for $k \geq l$, as these prior distributions are required to calculate $\tilde{z}_{l+1:L}$ and the discretization of the latent space. Also note that interval construction is *implicit*—we never have to do any computation over the set of all intervals, or store that set in memory, everything we need is computed efficiently only when required.

Since we discretize each latent layer into intervals of equal mass under the prior, the prior distribution over the indices i_l reduces to a uniform distribution over the interval indices, $U(i_l)$, which is not dependent on $i_{\neq l}$. This allows us to push/pop the i_l according to the prior in parallel. The full encoding and decoding procedures for a hierarchical latent model with the dynamic discretization are shown in Table B.1. Note that the operations in the two tables are ordered top to bottom.

Variables	Operation	Operation	Variables
x	$i_L \leftarrow Q_L(\cdot x)$	$i_{1:L} \leftarrow U(\cdot)$	$i_{1:L}$
x, i_L	$i_{L-1} \leftarrow Q_{L-1}(\cdot i_L, x)$	$x \leftarrow p(\cdot \tilde{z}_{1:L}(i_{1:L}))$	$x, i_{1:L}$
\vdots	\vdots	$i_1 \rightarrow Q_1(\cdot i_{2:L}, x)$	$x, i_{2:L}$
$x, i_{2:L}$	$i_1 \leftarrow Q_1(\cdot i_{2:L}, x)$	$i_2 \rightarrow Q_2(\cdot i_{3:L}, x)$	$x, i_{3:L}$
$x, i_{1:L}$	$x \rightarrow p(\cdot \tilde{z}_{1:L}(i_{1:L}))$	\vdots	\vdots
$i_{1:L}$	$i_{1:L} \rightarrow U(\cdot)$	$i_L \rightarrow Q_L(\cdot x)$	x
	(a) Encoding	(b) Decoding	

Table B.1: The BB-ANS encoding and decoding operations, in order from the top, for a hierarchical latent model with L layers. The Q_l are posterior distributions over the indices i_l of the discretized latent space for the l th latent, z_l . The discretization for the l th latent is created such that the intervals have equal mass under the prior.

Appendix C

The ResNet VAE architecture

A full description of the RVAE architecture is given in D. P. Kingma et al. (2016), and a full implementation can be found in our repository <https://github.com/hilloc-submission/hilloc>, but we give a short description below.

The RVAE is a hierarchical latent variable model, trained by maximization of the usual evidence lower bound (ELBO) on the log-likelihood:

$$\log p(x) \geq \mathbb{E}_{q(z|x)} \left[\log \frac{b(x, z)}{q(z|x)} \right]. \quad (\text{C.1})$$

We use L to denote the depth of the latent hierarchy, and label the latent layers $z_{1:L}$. There are skip connections in both the generative model, $p(x, z_{1:L})$, and the inference model, $q(z_{1:L}|x)$. Due to our requirement of using dynamic discretization, we use a top-down inference model¹. This means that we can factorize p and q

$$p(x, z_{1:L}) = p(x | z_{1:L}) p(z_L) \prod_{l=1}^{L-1} p(z_l | z_{l+1:L}) \quad (\text{C.2})$$

$$q(z_{1:L} | x) = q(z_L | x) \prod_{l=1}^{L-1} q(z_l | z_{l+1:L}, x) \quad (\text{C.3})$$

and express the ELBO as

$$\log p(x) \geq \mathbb{E}_{q(z_{1:L}|x)} [\log p(x | z_{1:L})] - D_{\text{KL}}(q(z_L | x) \| p(z_L)) \quad (\text{C.4})$$

$$- \sum_{l=1}^{L-1} \mathbb{E}_{q(z_{l+1:L}|x)} [D_{\text{KL}}(q(z_l | z_{l+1:L}, x) \| p(z_l | z_{l+1:L}))]. \quad (\text{C.5})$$

Where D_{KL} is the KL divergence. As in D. P. Kingma et al. (2016), the KL divergence terms are individually clamped by $\max(D_{\text{KL}}, \lambda)$, where λ is some

¹Note that in D. P. Kingma et al., 2016, this is referred to as ‘bidirectional inference’.

constant. This is an optimization technique known as *free bits*, and aims to prevent latent layers in the hierarchy collapsing to the prior.

Each layer in the hierarchy consists of a ResNet block with two sets of activations. One set of activations are calculated bottom-up (in the direction of x to z_L), and the other are calculated top-down. The bottom-up activations are used only within $q(z_{1:L} | x)$, whereas the top-down activations are used by both $q(z_{1:L} | x)$ and $p(x, z_{1:L})$. Every conditional distribution on a latent z_l is parameterized as a diagonal Gaussian distribution, with mean and covariance a function of the activations within the ResNet block, and the conditional distribution on x is parameterized by a discretized logistic distribution. Given activations for previous ResNet blocks, the activations at the following ResNet block are a combination of stochastic and deterministic features of the previous latent layer, as well as from skip connections directly passing the previous activations. The features are calculated by convolutions.

Note also that all latent layers are the same shape. Since we retained the default hyperparameters from the original implementation, each latent layer has 32 channels and spatial dimensions half those of the input (e.g. $\frac{h}{2} \times \frac{w}{2}$ for input of shape $h \times w$).

Bibliography

- Alakuijala, Jyrki et al. (2019). JPEG XL Next-Generation Image Compression Architecture and Coding Tools. In *Applications of Digital Image Processing XLII*. Vol. 11137, 111370K.
- Ballé, Johannes, Laparra, Valero, and Simoncelli, Eero P. (2017). End-to-End Optimized Image Compression. In *International Conference on Learning Representations (ICLR)*.
- Ballé, Johannes, Minnen, David, Singh, Saurabh, Hwang, Sung Jin, and Johnston, Nick (2018). Variational Image Compression with a Scale Hyperprior. In *International Conference on Learning Representations (ICLR)*.
- Benford, Frank (1938). The Law of Anomalous Numbers. In *Proceedings of the American Philosophical Society* 78.4, pp. 551–572.
- Bishop, Christopher (2006). *Pattern Recognition and Machine Learning*.
- Bloom, Charles (2014). Cbloom Rants: 02-10-14 - Understanding ANS - 9 (Blog Post).
- Bradbury, James, Frostig, Roy, Hawkins, Peter, Johnson, Matthew J., Leary, Chris, Maclaurin, Dougal, and Wanderman-Milne, Skye (2018). JAX: Composable Transformations of Python+NumPy Programs. Google.
- Brown, Tom B. et al. (2020). Language Models Are Few-Shot Learners. In *arXiv:2005.14165 [cs]*. arXiv: 2005.14165 [cs].
- Child, Rewon (2020). Very Deep VAEs Generalize Autoregressive Models and Can Outperform Them on Images. In *International Conference on Learning Representations*.
- Cover, Thomas. M. and Thomas, Joy A. (1991). *Elements of Information Theory*. Wiley Series in Telecommunications.

- Deng, Jia, Dong, Wei, Socher, Richard, Li, Li-Jia, Li, Kai, and Li Fei-Fei (2009). ImageNet: A Large-Scale Hierarchical Image Database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255.
- Duda, Jarek (2009). *Asymmetric Numeral Systems*. arXiv: 0902.0271 [cs, math].
- Duda, Jarek and Niemiec, Marcin (2016). *Lightweight Compression with Encryption Based on Asymmetric Numeral Systems*. arXiv: 1612.04662 [cs, math].
- Duda, Jarek, Tahboub, Khalid, Gadgil, Neeraj J., and Delp, Edward J. (2015). The Use of Asymmetric Numeral Systems as an Accurate Replacement for Huffman Coding. In *2015 Picture Coding Symposium (PCS)*, pp. 65–69.
- Frey, Brendan J. (1997). Bayesian Networks for Pattern Classification, Data Compression, and Channel Coding. PhD thesis. University of Toronto.
- Frey, Brendan J. and Hinton, Geoffrey E. (1996). Free Energy Coding. In *Proceedings of Data Compression Conference - DCC '96*, pp. 73–81.
- Germain, Mathieu, Gregor, Karol, Murray, Iain, and Larochelle, Hugo (2015). MADE: Masked Autoencoder for Distribution Estimation. In *International Conference on Machine Learning*. Chap. Machine Learning, pp. 881–889.
- Giesen, Fabian (2014). *Interleaved Entropy Coders*. arXiv: 1402.3392 [cs, math].
- Goodfellow, Ian, Bengio, Yoshua, and Courville, Aaron (2016). Deep Learning.
- Graves, Alex (2014). *Generating Sequences With Recurrent Neural Networks*. arXiv: 1308.0850.
- Gregor, Karol, Danihelka, Ivo, Graves, Alex, Rezende, Danilo, and Wierstra, Daan (2015). DRAW: A Recurrent Neural Network For Image Generation. In *International Conference on Machine Learning*, pp. 1462–1471.
- Han, Song, Mao, Huizi, and Dally, William J. (2016). Deep Compression: Compressing Deep Neural Network with Pruning, Trained Quantization and Huffman Coding. In *International Conference on Learning Representations (ICLR)*.

- Hinton, Geoffrey E. and van Camp, Drew (1993). Keeping Neural Networks Simple by Minimizing the Description Length of the Weights. In *Proceedings of the Sixth Annual Conference on Computational Learning Theory. COLT '93*, pp. 5–13.
- Ho, Jonathan, Lohn, Evan, and Abbeel, Pieter (2019). Compression with Flows via Local Bits-Back Coding. In *Advances in Neural Information Processing Systems 32*, pp. 3874–3883.
- Hoffman, Matthew and Johnson, Matthew J. (2016). ELBO Surgery: Yet Another Way to Carve up the Variational Evidence Lower Bound. In *Advances in Approximate Bayesian Inference (NeurIPS Workshop)*.
- Hoogeboom, Emiel, Peters, Jorn, van den Berg, Rianne, and Welling, Max (2019). Integer Discrete Flows and Lossless Compression. In *Advances in Neural Information Processing Systems 32*, pp. 12134–12144.
- Hubara, Itay, Courbariaux, Matthieu, Soudry, Daniel, El-Yaniv, Ran, and Bengio, Yoshua (2016). Binarized Neural Networks. In *Advances in Neural Information Processing Systems 29*, pp. 4107–4115.
- Huffman, David A. (1952). A Method for the Construction of Minimum-Redundancy Codes. In *Proceedings of the IRE* 40.9, pp. 1098–1101.
- Jain, Ajay, Abbeel, Pieter, and Pathak, Deepak (2020). Locally Masked Convolution for Autoregressive Models. In *Conference on Uncertainty in Artificial Intelligence*, pp. 1358–1367.
- Johnson, Matthew J., Duvenaud, David K., Wiltschko, Alex, Adams, Ryan P., and Datta, Sandeep R. (2016). Composing Graphical Models with Neural Networks for Structured Representations and Fast Inference. In *Advances in Neural Information Processing Systems 29*, pp. 2946–2954.
- Jun, Heewoo, Child, Rewon, Chen, Mark, Schulman, John, Ramesh, Aditya, Radford, Alec, and Sutskever, Ilya (2020). Distribution Augmentation for Generative Modeling. In *International Conference on Machine Learning*.
- Kingma, Diederik P., Salimans, Tim, Jozefowicz, Rafal, Chen, Xi, Sutskever, Ilya, and Welling, Max (2016). Improved Variational Inference with Inverse Autoregressive Flow. In *Advances in Neural Information Processing Systems 29*, pp. 4743–4751.

- Kingma, Diederik P. and Welling, Max (2014). Auto-Encoding Variational Bayes. In *International Conference on Learning Representations (ICLR)*.
- Kingma, Friso H., Abbeel, Pieter, and Ho, Jonathan (2019). Bit-Swap: Recursive Bits-Back Coding for Lossless Compression with Hierarchical Latent Variables. In *International Conference on Machine Learning*.
- Le Paine, Tom, Khorrami, Pooya, Chang, Shiyu, Zhang, Yang, Ramachandran, Prajit, Hasegawa-Johnson, Mark A., and Huang, Thomas S. (2016). *Fast Wavenet Generation Algorithm*. arXiv: 1611.09482 [cs].
- Lecun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-Based Learning Applied to Document Recognition. In *Proceedings of the IEEE* 86.11, pp. 2278–2324.
- Liu, Feng, Hernandez-Cabronero, Miguel, Sanchez, Victor, Marcellin, Michael W., and Bilgin, Ali (2017). The Current Role of Image Compression Standards in Medical Imaging. In *Information* 8.4, p. 131.
- Maaløe, Lars, Fraccaro, Marco, Liévin, Valentin, and Winther, Ole (2019). BIVA: A Very Deep Hierarchy of Latent Variables for Generative Modeling. In *Advances in Neural Information Processing Systems 32*, pp. 6551–6562.
- MacKay, David J. C. (2003). *Information Theory, Inference and Learning Algorithms*.
- Maclaurin, Dougal, Duvenaud, David, Johnson, Matthew J., and Townsend, James (2015). Autograd: Efficiently Compute Derivatives of Numpy Code. Harvard Intelligent Probabilistic Systems Group.
- McAnlis, Colton and Haecky, Aleks (2016). *Understanding Compression*.
- Menick, Jacob and Kalchbrenner, Nal (2018). Generating High Fidelity Images with Subscale Pixel Networks and Multidimensional Upscaling. In *International Conference on Learning Representations (ICLR)*.
- Mentzer, Fabian, Agustsson, Eirikur, Tschannen, Michael, Timofte, Radu, and Van Gool, Luc (2019). Practical Full Resolution Learned Lossless Image Compression. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10621–10630.

- Minnen, David, Ballé, Johannes, and Toderici, George D. (2018). Joint Autoregressive and Hierarchical Priors for Learned Image Compression. In *Advances in Neural Information Processing Systems 31*, pp. 10771–10780.
- Moffat, Alistair, Neal, Radford M., and Witten, Ian H. (1998). Arithmetic Coding Revisited. In *ACM Transactions on Information Systems (TOIS)* 16.3, pp. 256–294.
- Murphy, Kevin (2012). *Machine Learning: A Probabilistic Perspective*.
- Oliphant, Travis E. (2015). *Guide to NumPy*. Second.
- Ramachandran, Prajit, Le Paine, Tom, Khorrami, Pooya, Babaeizadeh, Mohammad, Chang, Shiyu, Zhang, Yang, Hasegawa-Johnson, Mark A., Campbell, Roy H., and Huang, Thomas S. (2017). *Fast Generation for Convolutional Autoregressive Models*. arXiv: 1704.06001 [cs, stat].
- Reed, Scott, van den Oord, Aaron, Kalchbrenner, Nal, Colmenarejo, Sergio G., Wang, Ziyu, Chen, Yutian, Belov, Dan, and Freitas, Nando (2017). Parallel Multiscale Autoregressive Density Estimation. In *International Conference on Machine Learning*. Chap. Machine Learning, pp. 2912–2921.
- Rezende, Danilo Jimenez, Mohamed, Shakir, and Wierstra, Daan (2014). Stochastic Backpropagation and Approximate Inference in Deep Generative Models. In *International Conference on Machine Learning*, pp. 1278–1286.
- Roelofs, Greg (1999). *PNG: The Definitive Guide*.
- Ruan, Yangjun, Ullrich, Karen, Severo, Daniel, Townsend, James, Khisti, Ashish, Doucet, Arnaud, Makhzani, Alireza, and Maddison, Chris J. (2021). *Improving Lossless Compression Rates via Monte Carlo Bits-Back Coding*. arXiv: 2102.11086 [cs, math, stat].
- Russakovsky, Olga et al. (2015). ImageNet Large Scale Visual Recognition Challenge. In *International Journal of Computer Vision* 115.3, pp. 211–252.
- Salakhutdinov, Ruslan and Murray, Iain (2008). On the Quantitative Analysis of Deep Belief Networks. In *International Conference on Machine Learning*.
- Salimans, Tim, Karpathy, Andrej, Chen, Xi, and Kingma, Diederik P. (2016). PixelCNN++: Improving the PixelCNN with Discretized Logistic Mix-

- ture Likelihood and Other Modifications. In *International Conference on Learning Representations (ICLR)*.
- Shannon, Claude. E. (1948). A Mathematical Theory of Communication. In *The Bell System Technical Journal* 27.3, pp. 379–423.
- Sneyers, Jon and Wuille, Pieter (2016). FLIF: Free Lossless Image Format Based on MANIAC Compression. In *2016 IEEE International Conference on Image Processing (ICIP)*, pp. 66–70.
- Sønderby, Casper K., Raiko, Tapani, Maaløe, Lars, Sønderby, Søren K., and Winther, Ole (2016). Ladder Variational Autoencoders. In *Advances in Neural Information Processing Systems 29*, pp. 3738–3746.
- Steinruecken, Christian (2014). Lossless Data Compression. PhD thesis. University of Cambridge.
- Theis, Lucas, Shi, Wenzhe, Cunningham, Andrew, and Huszár, Ferenc (2017). Lossy Image Compression with Compressive Autoencoders. In *International Conference on Learning Representations (ICLR)*.
- Townsend, James (2020). *A Tutorial on the Range Variant of Asymmetric Numeral Systems*. arXiv: 2001.09186 [cs, math, stat].
- Townsend, James, Bird, Thomas, and Barber, David (2019). Practical Lossless Compression with Latent Variables Using Bits Back Coding. In *International Conference on Learning Representations (ICLR)*.
- Townsend, James, Bird, Thomas, Kunze, Julius, and Barber, David (2020). HiLLoC: Lossless Image Compression with Hierarchical Latent Variable Models. In *International Conference on Learning Representations (ICLR)*.
- Townsend, James and Murray, Iain (2021). Lossless Compression with State Space Models Using Bits Back Coding. In *arXiv:2103.10150 [cs, math, stat]*. arXiv: 2103.10150 [cs, math, stat].
- Tschannen, Michael, Agustsson, Eirikur, and Lucic, Mario (2018). Deep Generative Models for Distribution-Preserving Lossy Compression. In *Advances in Neural Information Processing Systems 31*, pp. 5929–5940.

- Ullrich, Karen, Meeds, Edward, and Welling, Max (2017). Soft Weight-Sharing for Neural Network Compression. In *International Conference on Learning Representations (ICLR)*.
- Vahdat, Arash and Kautz, Jan (2020). NVAE: A Deep Hierarchical Variational Autoencoder. In *arXiv:2007.03898 [cs, stat]*. arXiv: 2007.03898 [cs, stat].
- van den Berg, Rianne, Gritsenko, Alexey A., Dehghani, Mostafa, Sønderby, Casper Kaae, and Salimans, Tim (2020). IDF++: Analyzing and Improving Integer Discrete Flows for Lossless Compression. In *arXiv:2006.12459 [cs, stat]*. arXiv: 2006.12459 [cs, stat].
- van den Oord, Aäron, Dieleman, Sander, Zen, Heiga, Simonyan, Karen, Vinyals, Oriol, Graves, Alex, Kalchbrenner, Nal, Senior, Andrew, and Kavukcuoglu, Koray (2016a). *WaveNet: A Generative Model for Raw Audio*. arXiv: 1609.03499.
- van den Oord, Aäron, Kalchbrenner, Nal, and Kavukcuoglu, Koray (2016b). Pixel Recurrent Neural Networks. In *International Conference on Machine Learning*, pp. 1747–1756.
- Voigtländer, Janis (2009). Bidirectionalization for Free! In *36th Symposium on Principles of Programming Languages, Savannah, Georgia, Proceedings*. Vol. 44. SIGPLAN Notices, pp. 165–176.
- Wallace, Chris S. (1990). Classification by Minimum-Message-Length Inference. In *Advances in Computing and Information*, pp. 72–81.
- Wallace, Gregory K. (1991). The JPEG Still Picture Compression Standard. In *Communications of the ACM* 34.4, pp. 30–44.
- Witten, Ian H., Neal, Radford M., and Cleary, John G. (1987). Arithmetic Coding for Data Compression. In *Communications of the ACM* 30.6, pp. 520–540.