

---

# Appendix - Compression with Bayesian Implicit Neural Representations

---

Anonymous Author(s)

Affiliation

Address

email

1 In addition to the four appendix sections mentioned in our main paper, we would like to draw atten-  
2 tion to two additional experiments: one evaluating the practical training and coding time, and the  
3 other investigating the impact of the number of training samples. These two experiments, especially  
4 the later one, offer **crucial** insights and are detailed in Appendix E1 and Appendix E2, respectively.

## 5 A Relative Entropy Coding with A\* Coding

---

### Algorithm 1 A\* encoding

---

**Require:** Proposal distribution  $p_{\mathbf{w}}$  and target distribution  $q_{\mathbf{w}}$ .

**Initialize :**  $N, G_0, \mathbf{w}^*, N^*, L \leftarrow 2^{|C|}, \infty, \perp, \perp, -\infty$

**for**  $i = 1, \dots, N$  **do**  $\triangleright N$  samples from proposal distribution

$\mathbf{w}_i \sim p_{\mathbf{w}}$

$G_i \sim \text{TruncGumbel}(G_{i-1})$

$L_i \leftarrow G_i + \log(q_{\mathbf{w}}(\mathbf{w}_i)/p_{\mathbf{w}}(\mathbf{w}_i))$   $\triangleright$  Perturbed importance weight

**if**  $L_i \leq L$  **then**

$L \leftarrow L_i$

$\mathbf{w}^*, N^* \leftarrow \mathbf{w}_i, i$

**end if**

**end for**

**return**  $\mathbf{w}^*, N^*$   $\triangleright$  Transmit the index  $N^*$

---

6 Recall that we would like to communicate a sample from the variational posterior distribution  $q_{\mathbf{w}}$   
7 using the proposal distribution  $p_{\mathbf{w}}$ . In our experiments, we used *global-bound depth-limited A\**  
8 *coding* to achieve this [1]. We describe the encoding procedure in Algorithm 1 and the decoding  
9 procedure in Algorithm 2. For brevity, we refer to this particular variant of the algorithm as *A\**  
10 *coding* for the rest of the appendix.

11 A\* coding is an importance sampler that draws  $N$  samples  $\mathbf{w}_1, \dots, \mathbf{w}_N \sim p_{\mathbf{w}}$  from the pro-  
12 posal distribution  $p_{\mathbf{w}}$ , where  $N$  is a parameter we pick. Then, it computes the importance weights  
13  $r(\mathbf{w}_n) = q_{\mathbf{w}}(\mathbf{w}_n)/p_{\mathbf{w}}(\mathbf{w}_n)$ , and sequentially perturbs them with truncated Gumbel<sup>1</sup> noise:

$$\tilde{r}_n = r(\mathbf{w}_n) + G_n, \quad G_n \sim \text{TruncGumbel}(G_{n-1}), \quad G_0 = \infty \quad (1)$$

14 Then, it can be shown that by setting

$$N^* = \arg \max_{n \in [1:N]} \tilde{r}_n, \quad (2)$$

---

<sup>1</sup>The PDF of a standard Gumbel random variable truncated to  $(-\infty, b)$  is given by  $\text{TruncGumbel}(x | b) = \mathbf{1}[x \leq b] \cdot \exp(-x - \exp(-x) + \exp(-b))$ .

---

**Algorithm 2** A\* decoding

---

**Simulate**  $\{w_i\} = \{w_1, \dots, w_N\}$   $\triangleright$  Simulate  $N$  samples from  $p_w$  with the shared seed  
**Receive**  $N^*$

**return**  $w^* \leftarrow w_{N^*}$   $\triangleright$  Receive the approximate posterior sample

---

15 we have that  $w_{N^*} \sim \tilde{q}_w$  is approximately distributed according to the target, i.e.  $\tilde{q}_w \approx q_w$ . More  
 16 precisely, we have the following result:

17 **Lemma A.1** (Bound on the total variation between  $\tilde{q}_w$  and  $q_w$  (Lemma D.1 in [2])). *Let us set the*  
 18 *number of proposal samples simulated by Algorithm 1 to  $N = 2^{D_{\text{KL}}[q_w \| p_w] + t}$  for some parameter*  
 19  *$t \geq 0$ . Let  $\tilde{q}_w$  denote the approximate distribution of the algorithm's output for this choice of  $N$ .*  
 20 *Then,*

$$D_{TV}(q_w, \tilde{q}_w) \leq 4\epsilon, \quad (3)$$

21 where

$$\epsilon = 2^{-t/4} + 2 \frac{\mathbb{P}_{Z \sim q_w}[\log_2 r(Z) \geq D_{\text{KL}}[Q \| P] + t]}{\mathbb{P}_{Z \sim q_w}[\log_2 r(Z) \geq D_{\text{KL}}[Q \| P] + t]}^{1/2}. \quad (4)$$

22 This result essentially tells us that we should draw at least around  $2^{D_{\text{KL}}[q_w \| p_w]}$  samples to ensure  
 23 low sample bias, and beyond this, the bias decreases exponentially quickly as  $t \rightarrow \infty$ . However,  
 24 note that the number of samples we need also increases exponentially quickly with  $t$ . In practice,  
 25 we observed that when  $D_{\text{KL}}[q_w \| p_w]$  is sufficiently large (around 16-20 bits), setting  $t = 0$  already  
 26 gave good results. To encode  $N^*$ , we built an empirical distribution over indices using our training  
 27 datasets and used it for entropy coding to find the optimal variable-length code for the index.

28 In short, on the encoder side,  $N$  random samples are obtained from the proposal distribution  $p_w$ ,  
 29 and we select the sample  $w_i$  and transmit its index  $N^*$  that has the greatest perturbed importance  
 30 weight. On the decoder side, those  $N$  random samples can be simulated with the same seed held by  
 31 the encoder. The decoder only needs to find the sample with the index  $N^*$ . Therefore, the decoding  
 32 process of our method is very fast. We also provide the specific coding time in Appendix E1.

## 33 B Closed-Form Solution for Updating Model Prior

34 In this section, we derive the analytic expressions for the prior parameter updates in our iterative  
 35 prior learning procedure when both the prior and the posterior are Gaussian distributions. Given a  
 36 set of training data  $\{D_i\} = \{D_1, D_2, \dots, D_M\}$ , we fit a variational distribution  $q_w^{(i)}$  to represent each  
 37 of the  $D_i$ s. To do this, we minimize the loss (abbreviated as  $L$  later)

$$\mathbb{E}_\beta(\theta_p, \{q_w^{(i)}\}) = \frac{1}{M} \sum_{i=1}^M L_\beta(D_i, q_w^{(i)}, p_{w; \rho}) \quad (5)$$

$$= \frac{1}{M} \sum_{i=1}^M \left\{ \mathbb{E}_{w \sim q_w} [\Delta(\mathbf{y}, f(\mathbf{x} | \mathbf{w}))] + \beta \cdot D_{\text{KL}}[q_w \| p_{w; \rho}] \right\}. \quad (6)$$

38 Now calculate the derivative w.r.t. the prior distribution parameter  $p_{w; \rho}$ ,

$$\frac{\partial L}{\partial \theta_p} = \frac{1}{M} \sum_{i=1}^M \frac{\partial D_{\text{KL}}[q_w \| p_{w; \rho}]}{\partial \theta_p} \quad (7)$$

39 Considering we choose factorized Gaussian as variational distributions, the KL divergence is

$$D_{\text{KL}}[q_w^{(i)} \| p_{w; \rho}] = D_{\text{KL}}[\mathcal{N}(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i)) \| \mathcal{N}(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i))] \quad (8)$$

$$= \frac{1}{2} \log \frac{\sigma_p}{\sigma_q^{(i)}} + \frac{\sigma_q^{(i)} + (\boldsymbol{\mu}_q^{(i)} - \boldsymbol{\mu}_p)^2}{\sigma_p} - \frac{1}{2} \quad (9)$$

40 To compute the analytical solution, let

$$\frac{\partial L}{\partial \theta_p} = \frac{1}{M} \sum_{i=1}^M \frac{\partial D_{\text{KL}}[q_{\mathbf{w}} \| p_{\mathbf{w}, \rho}]}{\partial \theta_p} = 0. \quad (10)$$

41 Note here  $\sigma$  refers to variance rather than standard deviation. The above equation is equivalent to

$$\begin{aligned} \frac{\partial L}{\partial \mu_p} &= \sum_{i=1}^M \frac{\mu_p - \mu_q^{(i)}}{\sigma_p} = 0, \\ \frac{\partial L}{\partial \sigma_p} &= \sum_{i=1}^M \left[ \frac{1}{\sigma_p} - \frac{\sigma_q^{(i)} + (\mu_q^{(i)} - \mu_p)^2}{\sigma_p^2} \right] = 0. \end{aligned} \quad (11)$$

42 We finally can solve these equations and get

$$\mu_p = \frac{1}{M} \sum_{i=1}^M \mu_q^{(i)}, \quad \sigma_p = \frac{1}{M} \sum_{i=1}^M [\sigma_q^{(i)} + (\mu_q^{(i)} - \mu_p)^2] \quad (12)$$

43 as the result of Equation (5) in our main text. In short, this closed-form solution provides an efficient  
44 way to update the model prior from a bunch of variational posteriors. It makes our method simple  
45 in practice, unlike some previous methods [3, 4] that require expensive meta-learning loops.

## 46 C Dynamic Adjustment of $\beta$

47 When learning the model prior, the value of  $\beta$  that controlling the rate-distortion trade-off is defined  
48 in advance to train the model prior at a specific bitrate point. After obtaining the model prior, we  
49 will first partition the network parameters into  $K$  groups  $\mathbf{w}_{1:K} = \{\mathbf{w}_1, \dots, \mathbf{w}_K\}$  according to the  
50 average approximate coding cost of training data, as described in Section 3.3 of the main text. Now  
51 for training the variational posterior for a given test datum, to ensure the coding cost of each group  
52 is close to  $\kappa = 16$  bits, we adjust the value of  $\beta$  dynamically when optimizing the posteriors. The  
53 detailed algorithm is illustrated here in Algorithm 3.

---

### Algorithm 3 Dynamic $\beta$ adjustment for optimizing the posteriors

---

**Require:**  $\beta, \mathbf{w}_{1:K} = \{\mathbf{w}_1, \dots, \mathbf{w}_K\}$

**Initialize:**  $\lambda_k = \beta, k = 1, \dots, K$

**Initialize:** variational posterior  $q_{\mathbf{w}_k}, k = 1, \dots, K$

**for**  $i \leftarrow \text{NumberIter}$  **do**

$\delta_k = D_{\text{KL}}[q_{\mathbf{w}_k} \| p_{\mathbf{w}_k}], k = 1, \dots, K$

$q_{\mathbf{w}_{1:K}} \leftarrow \text{VariationalUpdate}(L_{\lambda_{1:K}}) \quad \triangleright L_{\lambda_{1:K}}$  is defined in Equation 8 in the main text

**if**  $(i \bmod 15) = 0$  **then**

**if**  $\delta_k > \kappa$  **then**  $\lambda_k = \lambda_k \cdot 1.05$

**end if**

**if**  $\delta_k < \kappa - 0.4$  **then**  $\lambda_k = \lambda_k / 1.05$

**end if**

**end if**

**end for**

**return**  $q_{\mathbf{w}_k}, \lambda_k, k = 1, \dots, K$

---

54 The algorithm is improved from Havasi et al. [5] to stabilize training, in the way that we set an  
55 interval  $[\kappa - 0.4, \kappa]$  as buffer area where we do not change the value of  $\lambda_k$ . Here we only adjust  $\lambda_k$   
56 every 15 iterations to avoid frequent changes at the initial training stage.

## 57 D Experiment Details

58 We introduce the experimental settings here and summarize the settings in Table 1.

	CIFAR-10	Kodak		LibriSpeech
		Smaller Model	Larger Model	
Network Structure				
number of MLP layer	4	6	7	6
hidden unit	16	48	56	48
Fourier embedding	32	64	96	64
number of parameters	1123	12675	21563	12675
Learning Model Prior from Training Data				
number of training data	2048	512	512	1024
epoch number	128	96	96	128
learning rate	0.0002	0.0001	0.0001	0.0002
iteration / epoch (except the first epoch)	100	200	200	100
iteration number in the first epoch	250	500	500	250
initialization of posterior variance	$9 \times 10^{-6}$	$4 \times 10^{-6}$	$4 \times 10^{-6}, 4 \times 10^{-10}$	$4 \times 10^{-9}$
$\beta$	$2 \times 10^{-5}, 5 \times 10^{-6}, 2 \times 10^{-6}$ $1 \times 10^{-6}, 5 \times 10^{-7}$	$10^{-7}, 10^{-8}, 4 \times 10^{-8}$	$4 \times 10^{-6}$	$10^{-7}, 3 \times 10^{-8}$ $10^{-8}, 10^{-9}$
Optimize the Posterior of a Test Datum				
iteration number	25000	25000	25000	25000
learning rate	0.0002	0.0001	0.0001	0.0002
training with 1/4 the points (pixels)	$\times$	$\checkmark$	$\checkmark$	$\times$
number of group (KL budget = 16 bits / group)	(58, 89, 146, 224, 285)	(1729, 2962, 3264)	(5503, 7176)	(1005, 2924, 4575, 6289)
bitrate, (bpp for images, Kbps for audios)	(0.91, 1.39, 2.28, 3.50, 4.45)	(0.070, 0.110, 0.132)	(0.224, 0.293)	(5.36, 15.59, 24.40, 33.54)
PSNR, dB	(0.91, 1.39, 2.28, 3.50, 4.45)	(0.070, 0.110, 0.132)	(0.224, 0.293)	(5.36, 15.59, 24.40, 33.54)

Table 1: Hyper parameters in our experiments.

## 59 D.1 CIFAR-10

60 We use a 4-layer MLP with 16 hidden units and 32 Fourier embeddings for the CIFAR-10 dataset.  
61 The model prior is trained with 128 epochs to ensure convergence. Here, the term “epoch” is used  
62 to refer to optimizing the posteriors and updating the prior in the Algorithm 1 in the main text. For  
63 each epoch, the posteriors of all 2048 training data are optimized for 100 iterations using the local  
64 reparameterization trick [6], except the first epoch that contains 250 iterations. We use the Adam  
65 optimizer with learning rate 0.0002. The posterior variances are initialized as  $9 \times 10^{-6}$ .

66 After obtaining the model prior, given a specific test CIFAR-10 image to be compressed, the pos-  
67 terior of this image is optimized for 25000 iterations, with the same optimizer. When we finally  
68 progressively compress and finetune the posterior, the posteriors of the uncompressed parameter  
69 groups are finetuned for 15 iterations with the same optimizer once a previous group is compressed.

## 70 D.2 Kodak

71 For Kodak dataset, since training on high-resolution image takes much longer time, the model prior  
72 is learned using fewer training data, i.e., only 512 cropped CLIC images [7]. We also reduce the  
73 learning rate of the Adam optimizer to 0.0001 to stabilize training. In each epoch, the posterior of  
74 each image is trained for 200 iterations, except the first epoch that contains 500 iterations. We also  
75 reduce the total epoch number to 96 which is empirically enough to learn the model prior.

76 We use two models with different capacity for compressing high-resolution Kodak images. The  
77 smaller model is a 6-layer SIREN with 48 hidden units and 64 Fourier embeddings. This model is

78 used to get the three low-bitrate points in Figure 2b in our main text, where the corresponding beta  
 79 is set as  $\{10^{-7}, 10^{-8}, 4 \times 10^{-8}\}$ . Another larger model comprises a 7-layer MLP with 56 hidden  
 80 units and 96 Fourier embeddings, which is used for evaluation at the two relatively higher bitrate  
 81 points in Figure 2b in our main text. The betas of these two models have the same value  $2 \times 10^{-9}$ .  
 82 We empirically adjust the variance initialization from the set  $\{4 \times 10^{-6}, 4 \times 10^{-10}\}$  and find they  
 83 can affect the converged bitrate and achieve good performance. In particular, the posterior variance  
 84 is initialized as  $4 \times 10^{-10}$  to reach the highest bitrate point in the rate-distortion curve. The posterior  
 85 variance of other bitrate-points on Kodak dataset are all initialized as  $4 \times 10^{-6}$ .

86 **Important note:** It required significant empirical effort to find the optimal parameter settings we  
 87 described above, hence our note in the Conclusion and Limitations section that Bayesian neural  
 88 networks are inherently sensitive to initialization [8].

### 89 D.3 LibriSpeech

90 We randomly crop 1024 audio samples from LibriSpeech “train-clean-100” set [9] for learning the  
 91 model prior and randomly crop 24 test samples from “test-clean” set for evaluation. The model  
 92 structure is the same as the small model used for compressing Kodak images. We evaluate on four  
 93 bitrate points by setting  $\beta = \{10^{-7}, 3 \times 10^{-8}, 10^{-8}, 10^{-9}\}$ . There are 128 epochs, and each epoch  
 94 has 100 iterations with learning rate as 0.0002. The first epoch has 250 iterations. In addition, the  
 95 posterior variance is initialized as  $4 \times 10^{-9}$ . The settings for optimizing and finetuning posterior of  
 96 a test datum are the same as the experiments on Kodak dataset.

## 97 E Supplementary Experimental Results

### 98 E.1 Evaluation of Training and Practical Coding Time

99 **Training Time.** The most time-consuming part of setting up COMBINER for a data modality is  
 100 running the iterative prior learning algorithm we propose in the main text. Furthermore, optimizing  
 101 the variational posteriors takes up the bulk of the learning process since updating the prior parameters  
 102 given the variational posteriors can be done efficiently using the formulae we derive in Appendix B.  
 103 However, such posterior optimization can be done in parallel, especially given that our INRs have  
 104 very few parameters. To train the model prior on the CIFAR-10 dataset, we can train the posteriors  
 105 of 2048 images together in a single V100 GPU. It only takes around **20 minutes** to train for 128  
 106 epochs with almost 100 iterations per epoch. For training the model prior with 512 cropped CLIC  
 107 images, due to the limit of GPU memory, we run multiple processes simultaneously on 4 GTX  
 108 1080 GPUs, where each process runs for a single image. The entire training time on CLIC dataset  
 109 consumes around 30 hours. We note that the training time on CLIC dataset could be significantly  
 110 reduced with additional engineering effort, but we have not had the opportunity to do so due to time  
 111 constraints.

112 **Coding Time.** To compress a test datum, we first optimize its INR’s variational posterior for 25,000  
 113 iterations. Such optimization process should also be included as a part of encoding time, similar  
 114 to COIN [10]. In addition, the progressive posterior refinement process also takes a long time.  
 115 Therefore, the encoding time of our method is very long. Note that the encoding time of relative  
 116 entropy coding is negligible compared with the optimization process because our model is very  
 117 small, and there are not so many parameter groups. As a result, we are able to evaluate all the  
 118 10,000 images from the CIFAR-10 test set in parallel using a CPU cluster. To decode the network  
 119 parameters, the decoder only needs to search for the sample according to the received index, which  
 120 is very fast. The practical encoding and decoding time is shown in Table 2.

	CIFAR, bpp = 0.91	CIFAR, bpp = 4.45	Kodak, bpp = 0.070	Kodak, bpp = 0.293
encoding time	~10 minutes	~20 minutes	~2 hours	~4 hours
decoding time	0.051 second	0.075 second	0.410 second	0.542 second

Table 2: Practical encoding and decoding time of a specific image on 1080Ti GPU.

121 We show both the encoding and decoding time of our method on different datasets at different  
 122 bitrates. In fact, the decoding time is mainly consumed for relative entropy decoding and inference

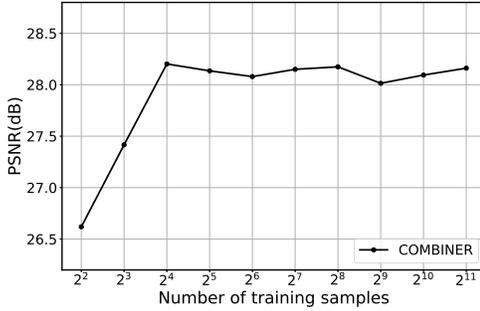


Figure 1: Impact of the number of training data.

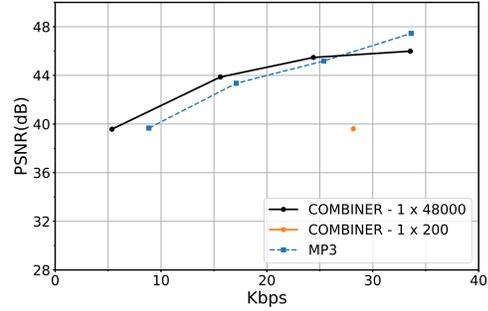


Figure 2: Compressing audios.

123 of the received MLP network. Usually, if there are more parameter groups, the coding time will be  
 124 longer. Therefore, decoding a Kodak image at our highest bitrate (0.293 bpp) consumes the most  
 125 decoding time (0.542 second), but is still very fast.

## 126 E.2 Number of Training Samples

127 Since the model prior is learned from a few training data, the number of training data may influence  
 128 the quality of the learned model prior. We train the model prior with a different number of training  
 129 images from the CIFAR-10 training set and evaluate the performance on 100 randomly selected test  
 130 images from the CIFAR-10 test set. Surprisingly, as shown in Figure 1, we found that even merely  
 131 16 training images can help to learn a good model prior. Considering the randomness of training  
 132 and testing, the performance on this test subset is almost the same when the number of training data  
 133 exceeds 16. This demonstrates that the model prior is quite robust and generalizes well to test data.  
 134 In our final experiments, the number of training samples is set to 2048 (on CIFAR-10 dataset) to  
 135 ensure the prior converges to a good optimum.

## 136 E.3 Compressing Audios with Small Chunks

137 The proposed approach does not need to compute the second-order gradient during training, which  
 138 helps to learn the model prior of the entire datum. Hence, compression with a single Bayesian INR  
 139 network helps to fully capture the global dependencies of a datum. That is the reason for our strong  
 140 performance on Kodak and LibriSpeech datasets. Here, we also conduct a group of experiment to  
 141 compare the influence of cropping audios into chunks. Unlike the experimental setting in our main  
 142 text that compresses every 3-second audio ( $1 \times 48000$ ) with a single MLP network, here we try to  
 143 crop all the 24 audios into small chunks, each of the chunk has the shape of  $1 \times 200$ . We use the same  
 144 network used for compressing CIFAR-10 images for our experiments here. As shown in Figure 2,  
 145 if we do not compress the audio as an entire entity, the performance will drops for around 5 dB.  
 146 It demonstrates the importance of compressing with a single MLP network to capture the inherent  
 147 redundancies within the entire data.

## 148 E.4 Additional Figures

149 We provide some examples of the decoded Kodak images in Figure 3.

