

Neural Distributed Image Compression with Cross-Attention Feature Alignment

Nitish Mital^{1,*}, Ezgi Özyılkan^{1,†}, Ali Garjani^{1,‡}, and Deniz Gündüz^{*}

^{*}Dept. of Electrical and Electronics Engineering, Imperial College London

[†]Dept. of Electrical and Computer Engineering, New York University

[‡]Section of Mathematics, EPFL

{n.mital, d.gunduz}@imperial.ac.uk, eo2135@nyu.edu, ali.garjani@epfl.ch

Abstract

We consider the problem of compressing an information source when a correlated one is available as side information only at the decoder side, which is a special case of the distributed source coding problem in information theory. In particular, we consider a pair of stereo images, which have overlapping fields of view, and are captured by a synchronized and calibrated pair of cameras as correlated image sources. In previously proposed methods, the encoder transforms the input image to a latent representation using a deep neural network, and compresses the quantized latent representation losslessly using entropy coding. The decoder decodes the entropy-coded quantized latent representation, and reconstructs the input image using this representation and the available side information. In the proposed method, the decoder employs a cross-attention module to align the feature maps obtained from the received latent representation of the input image and a latent representation of the side information. We argue that aligning the correlated patches in the feature maps allows better utilization of the side information. We empirically demonstrate the competitiveness of the proposed algorithm on KITTI and Cityscape datasets of stereo image pairs. Our experimental results show that the proposed architecture is able to exploit the decoder-only side information in a more efficient manner compared to previous works.

1. Introduction

Image compression is a fundamental task in image processing that aims to preserve the visual image content while reducing the bit rate needed for storage or transmission. The compression may be lossless, that is, when multiple

samples of the information source are compressed jointly such that the source can be reconstructed with a vanishing probability of error, or lossy, that is, allowing a non-zero distortion in the reconstruction in order to achieve higher compression rates. Shannon showed that the *entropy* of the source is a fundamental bound on the bit rate for lossless compression. In the lossy case, continuous-valued data (such as vectors of image pixel intensities) must be first quantized to a finite set of discrete values, which inherently introduces some degree of error. Therefore, for lossy compression, one must trade-off between two competing costs: the entropy of the discretized latent representation (rate) and the error arising from the quantization step (distortion). Traditional image compression schemes, like JPEG2000 [30] and BPG [6], typically consist of partitioning the image into small pre-determined blocks, which are processed through linear transforms like the discrete wavelet transform (DWT), in order to decorrelate the pixel values, and to obtain a latent representation of the image, followed by intra block prediction (motion search) and residual coding to exploit repetition and self-similarity of the image content, which reduces the entropy of its representation. This is then followed by quantizing the latent representation, and an entropy coder to store/send the resulting quantized representation most efficiently. On the other hand, recently proposed machine-learning-driven compression algorithms [3, 5, 16, 23, 26, 32, 33], which employ deep neural networks (DNNs), achieve impressive performance results, outperforming classical and standard methods, by decorrelating the image values with a nonlinear transform, parameterized by a DNN, in order to obtain a latent representation, which is then quantized and entropy coded using a learned probability distribution.

In this work, we are interested in DNN-aided *distributed stereo image compression*, where an image \mathbf{y} from the stereo image pair (\mathbf{x}, \mathbf{y}) is available as side information only at

¹Contributed equally to this work.

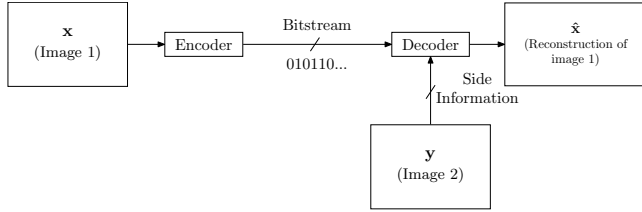


Figure 1: System model.

the decoder side (see Fig. 1). This scenario can occur, for example, when there are multiple distributed unmanned aerial vehicles, autonomous vehicles, or simply multiple static cameras that capture images with overlapping fields of view. Note that these captured images are highly correlated due to overlapping fields of view. Assume that one of the cameras delivers its image (in a lossless fashion) to the destination, e.g., a central storage or a processing unit. The other camera, instead of employing a standard single image compression algorithm, should be able to benefit from the presence of a highly correlated image obtained from the first camera, even though it does not have direct access to this side information image at the encoder side. This is a special case of the more general distributed source coding (DSC) problem where two distributed encoders communicate their sources to the decoder, characterized by an achievable region of rate pairs (R_x, R_y) $((D_x, D_y)$, where D_x and D_y denote the distortions in the reconstruction of sources \mathbf{x} and \mathbf{y} , respectively. The case we consider here is a corner point of the achievable rate region corresponding to $D_y = 0$ (lossless compression), implying that $R_y = H(\mathbf{y})$. The benefit of decoder-only side information in compression was first characterized by Slepian and Wolf in [31] for the lossless compression case, and by Wyner and Ziv in [39] for the lossy compression case.

1.1. Related work

1.1.1 Single image compression

There has been a surge of interest in DNN models for image compression, most notably the ones proposed in [3, 5, 16, 23, 26, 32, 33]. In [3], an autoencoder-based model with a parameterized distribution as prior for the latent is trained with a rate-distortion loss for a fixed target bit rate. An extension is proposed in [5] that introduces a hyperprior to capture the spatial dependencies between the elements of the latent representation by estimating their standard deviations, thus enabling better compression of the latent representation by the entropy coder. In [16], context-adaptive entropy models are introduced, while in [23], an autoregressive network is used as a non-factorized conditional entropy model. Both [16, 23] generalize the hyperpriors of [5] to estimate both the mean and variance of the Gaussian priors of the latent representation. Many other

approaches and architectures have been recently proposed, like saliency-driven compression [26], dense-blocks and content-weighting [18], non-local attention [7, 8, 15], and generative-adversarial networks (GANs) [1, 20].

1.1.2 Attention

Self-attention in vision applications was first introduced in [12], where the image is tiled into a sequence of flattened patches, and attention is applied to this sequence of patches. In [27], attention mechanism is restricted to a local neighborhood in order to fully replace convolutional layers. For rectified stereo images, in particular, stereo attention modules (SAM) are introduced in [40] for stereo image super-resolution, where attention at a certain location in the left (or right) image is limited to the corresponding epipolar line in the other image.

1.1.3 Centralized stereo compression

Centralized stereo image compression was first considered in the DSIC model in [19], and the HESIC model in [10], in which both the left and right images are available at the encoder, and are jointly compressed. In DSIC, a dense warp field is estimated using disparity estimation between the two images, and warped features from the left image are fed to the encoder and decoder of the right image. In HESIC, the right image is warped by an estimated homography, and only its residual with respect to the first image is encoded. Subsequent works include the SASIC model [38], and the bi-directional contextual transform module (Bi-CTM) and a bi-directional conditional entropy model (Bi-CEM) [17]. The SASIC approach computes the optimal horizontal shifts of each channel of the latent representation to match the second image, and then encodes only the residual of the shifted channel with respect to the corresponding channel of the second image. SASIC also connects the encoder-decoder pipelines of the two images with stereo attention modules [40]. In [17], the main idea is to avoid the limitation of sequential coding of the two stereo images by introducing an “inter-view context dependency” mechanism.

1.1.4 Distributed stereo compression

In the current literature, the DNN-based methods that explicitly address distributed (stereo) image compression are: (1) the DRASIC model in [11], (2) the DSIN model in [2], (3) the NDIC model in [24], and (4) vector quantized variational autoencoder-based approach in [36]. In [2], the authors exploit high spatial correlations between pairs of stereo images, having significantly overlapping fields of view. By finding corresponding patches between an intermediate reconstructed image and the side information im-

age, and computing their correlations, the authors then use these patches to refine the reconstructed image at the decoder. The process of finding the corresponding patches is non-differentiable, since it is done by using the $\text{argmax}(\cdot)$ function, which possibly prevents the network from learning the inter-dependencies between the images in an optimal way.

The paper [24] uses a different approach by explicitly modeling the correlation between the two stereo images. More precisely, [24] models the two images as being generated by a common set of features, as well as two independent sets of features that capture the information in the respective images that is not captured by the set of common features. In order to minimize the redundant information that is transmitted, the encoder only sends the independent information corresponding to the input image, while the set of common features between the input image and the side information are recovered locally only from the latter. The paper [36] uses a vector quantized variational autoencoder (VQ-VAE) where, unlike most existing DNN-based image compression schemes that use *uniform quantization*, the model learns the quantization codebook, that is, it employs *non-uniform quantization*. In [11], the authors propose a framework for distributed compression of correlated sources, followed by joint decoding, by employing a recurrent autoencoder architecture that processes the residual content over repeated multiple iterations in order to achieve better reconstruction performance.

2. Proposed method

2.1. Main contribution

In this paper, we use the NDIC model [24] as the “backbone”, which itself is built upon the model in [3] as its backbone. In principle, any other single image compression algorithm can also be used as the backbone for NDIC; and hence, also for our model. In addition, we augment this backbone by introducing some transformer-based blocks, specifically *cross-attention modules* (CAMs) between the intermediate latent representations in different stages of the decoders of the input image and the side information, whose purpose is to align the corresponding patches. This is similar to the “patch-matching” idea proposed in [2], but our method provides a differentiable alternative to the search-based algorithm used in [2]. Unlike the SAM approach [40], used also in SASIC [38], the CAM technique we introduce computes the attention globally, between patches of the latent representations over all channels, similarly to [12, 34]. We show that our method outperforms the solution provided in [24]. We also show that our method is able to perform well in the case of unsynchronized and uncalibrated stereo cameras, that is, when the correlated images are generated at different time steps.

2.2. Architecture

In this section, we describe the autoencoder architecture we use in our compression scheme. Following the method proposed in [24], we model the images \mathbf{x} and \mathbf{y} as being generated by random variables \mathbf{w} , \mathbf{v}_x and \mathbf{v}_y . The variable \mathbf{w} is meant to capture the common features between the two images, while the variables \mathbf{v}_x and \mathbf{v}_y , which are called the private information variables of the respective images \mathbf{x} and \mathbf{y} , are designated to capture the private aspects of \mathbf{x} and \mathbf{y} that are *not* captured by the common variable \mathbf{w} . The decoder reconstructs not only the required image, but also the side information image from \mathbf{w} and \mathbf{v}_y , in order to ensure that the common features \mathbf{w} , extracted from \mathbf{y} only, are relevant to both images. The common information \mathbf{w} here is defined in the sense of Witsenhausen, Gacs and Korner [13, 37], where it corresponds to a deterministic function $\mathbf{w} = f(\mathbf{y})$ ($= f'(\mathbf{x})$) of the two information sources, that is, two separate observers of \mathbf{x} and \mathbf{y} are able to agree as to the value of \mathbf{w} with probability one.

See Fig. 2 for an illustration of the proposed distributed compression algorithm using CAMs. The encoder maps the image \mathbf{x} to a latent representation \mathbf{v}_x by applying a transform \mathbf{g}_{ax} , which is parameterized by weights ϕ_x . Then, the latent representation \mathbf{v}_x is quantized to obtain $\hat{\mathbf{v}}_x \in \mathbb{Z}^m$, where its elements are rounded to the closest integer values. Since the quantization step is a non-differentiable operation, which prevents end-to-end training, it is instead replaced by additive uniform random noise over $[-0.5, 0.5]$ during training (see [3] for a similar reasoning). Thus, \mathbf{v}_x is perturbed by uniform noise during training to obtain $\tilde{\mathbf{v}}_x$, which approximates the quantized latents $\hat{\mathbf{v}}_x$. Similarly to [24], the decoder extracts $\mathbf{w} = \mathbf{f}(\mathbf{y}; \phi_f)$ by applying a nonlinear transform \mathbf{f} to image \mathbf{y} , where ϕ_f refers to the weights of the respective DNN. During training, the transform \mathbf{f} learns to extract features from the SI that estimate the common information between the stereo images. At the decoder, \mathbf{w} is concatenated with the received latent variable $\tilde{\mathbf{v}}_x$, and is given as an input to the first layer of the primary image’s decoder network \mathbf{g}_{sx} , denoted by $\mathbf{g}_{sx}^{(1)}$, which is parameterized by weights $\theta_x^{(1)}$. Simultaneously, the side information image’s decoder maps the correlated image \mathbf{y} to the latent representation \mathbf{v}_y using a transform \mathbf{g}_{ay} , which is parameterized by weights ϕ_y . It then concatenates the common variable \mathbf{w} with \mathbf{v}_y , and then inputs it to the first layer of a decoder network \mathbf{g}_{sy} , denoted by $\mathbf{g}_{sy}^{(1)}$, which is parameterized by weights $\theta_y^{(1)}$.

In order to overcome the limitation of the convolutional layers to allow only local feature interaction between the two images, we introduce CAMs between the decoder pipelines of the two images, that capture global correlations between the intermediate latent representations in the decoder architectures for both images. Then the outputs from

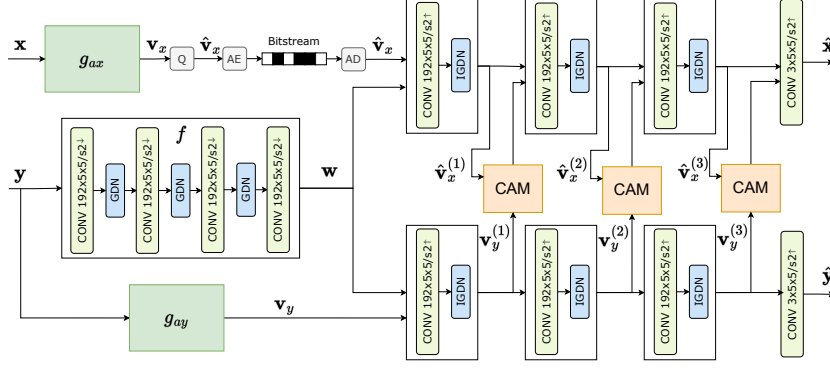


Figure 2: Proposed model architecture.

$\mathbf{g}_{sx}^{(1)}$ and $\mathbf{g}_{sy}^{(1)}$ are fed as inputs into a CAM (described in detail at Section 2.3), that morphs and aligns the output feature map from $\mathbf{g}_{sy}^{(1)}$, denoted by $\mathbf{v}_y^{(1)}$, with the output feature map obtained from $\mathbf{g}_{sx}^{(1)}$, denoted by $\hat{\mathbf{v}}_x^{(1)}$. Next, the output of the CAM, i.e., $\mathbf{v}_{CAM}^{(1)}$, is concatenated with $\hat{\mathbf{v}}_x^{(1)}$, and fed to the second layer $\mathbf{g}_{sx}^{(2)}$. As seen in Fig. 2, this procedure is repeated in the next two consecutive layers. In general, the outputs of the i^{th} layer of the decoder networks, which are $\hat{\mathbf{v}}_x^{(i)}$ and $\mathbf{v}_y^{(i)}$, are fed into a CAM, whose output, i.e., $\mathbf{v}_{CAM}^{(i)}$, is concatenated with $\hat{\mathbf{v}}_x^{(i)}$ in order to be fed to the $(i+1)^{th}$ layer of \mathbf{g}_{sx} . The reconstructed input image $\hat{\mathbf{x}}$ and the reconstructed side information image $\hat{\mathbf{y}}$ are obtained as the outputs of the decoder blocks \mathbf{g}_{sx} and \mathbf{g}_{sy} , respectively. Note that the latent representation \mathbf{v}_y is neither quantized nor perturbed with uniform noise, unlike \mathbf{v}_x . This is because the encoding and decoding of image \mathbf{y} happen at the decoder side without it being transmitted over the channel. During training, we minimize the following loss function

$$L = R_x + \lambda D_x + \alpha(R_y + \lambda D_y) + \beta R_w, \quad (1)$$

where R_x, R_y and R_w are the entropy estimates of $\mathbf{v}_x, \mathbf{v}_y$ and \mathbf{w} , respectively, and D_x and D_y are the distortion terms for the reconstructions of the input image and the side information, respectively. In particular, R_x represents the rate of transmission of the input image \mathbf{x} . Similarly to previous works [3, 24], the probability distributions of the variables \mathbf{w}, \mathbf{v}_x and \mathbf{v}_y are modeled using univariate non-parametric, fully factorized density functions, which are used to compute the associated entropy terms. In Eq. (1), the hyperparameter β controls how much importance is given to the complexity of the common information to be extracted by the decoder, and α determines how much emphasis is given to the reconstruction loss of the side information. Since our main objective is the reconstruction of only \mathbf{x} , we argue that the terms $R_y + \lambda D_y$ and R_w act as regularizers for the main objective under consideration, that is the rate-distortion per-

formance of the primary image \mathbf{x} .

2.3. Cross-attention module (CAM)

The CAM takes as input the tensors $\mathbf{v}_x^{(i)}, \mathbf{v}_y^{(i)} \in \mathbb{R}^{C \times H \times W}$, where C is the number of channels, H is the height, and W is the width. The input tensors are tiled into $N = \frac{CHW}{C_p H_p W_p}$ 3D patches of dimension $C_p \times H_p \times W_p$, where C_p is the number of channels, H_p is the height, and W_p is the width of each patch. Using a linear layer, the set of patches is transformed to a set of patch embeddings, denoted by $\mathbf{P}_x = (\mathbf{p}_x^1, \dots, \mathbf{p}_x^N) \in \mathbb{R}^{D_1 \times N}$ and $\mathbf{P}_y = (\mathbf{p}_y^1, \dots, \mathbf{p}_y^N) \in \mathbb{R}^{D_1 \times N}$ of $\mathbf{v}_x^{(i)}$ and $\mathbf{v}_y^{(i)}$, respectively, where D_1 is the length of each patch embedding. We define three learnable weight matrices, namely *query* ($\mathbf{W}_x^Q \in \mathbb{R}^{D_1 \times D_2}$), *key* ($\mathbf{W}_y^K \in \mathbb{R}^{D_1 \times D_2}$), and *value* ($\mathbf{W}_y^V \in \mathbb{R}^{D_1 \times D_2}$), where D_2 is the length of the query, key and value corresponding to each patch embedding. The patch embeddings are projected onto these weight matrices to obtain $\mathbf{Q}_x = (\mathbf{P}_x)^T \mathbf{W}_x^Q$, $\mathbf{K}_y = (\mathbf{P}_y)^T \mathbf{W}_y^K$, and $\mathbf{V}_y = (\mathbf{P}_y)^T \mathbf{W}_y^V$. Finally, the output of the CAM is computed as

$$\mathbf{v}_{CAM}^{(i)} = \text{Unpack embedding} \left(\text{Softmax} \left(\frac{\mathbf{Q}_x \mathbf{K}_y^T}{\sqrt{D_2}} \right) \mathbf{V}_y \right), \quad (2)$$

where the “unpack embedding” operation reverses the embedding operation done on the patches. See Fig. 3 for an overall summary of the CAM architecture. In the code, we employ a multi-headed attention mechanism, that is, multiple attention weights are computed in parallel, similarly to [34].

3. Experiments

3.1. Experimental setup

In order to assess the rate-distortion performance of our proposed approach with respect to the existing models for DSC as well as to the point-to-point neural compression

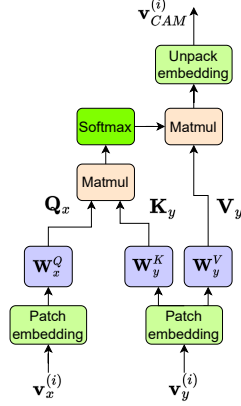


Figure 3: Cross-attention module architecture. The Matmul operation refers to matrix multiplication (see Eq. (2)).

baselines [3, 5], we conducted a number of experiments using the PyTorch framework [25]. Our code is publicly available¹.

See Fig. 2 for the proposed DNN architecture. The non-linear transforms \mathbf{g}_{ax} and \mathbf{g}_{ay} have the same structure as those proposed in [3]. More specifically, these transforms are consisting of convolutional layers followed by either linear (i.e., rectified linear unit) or nonlinear functions (i.e., generalized divisive normalization [GDN] [4] and inverse generalized divisive normalization [IGDN]). In [3], it has been shown that GDN and IGDN are particularly suited for density modelling within the context of neural image compression. Additionally, we introduce the transform denoted by \mathbf{f} , as proposed in [24], as well as CAMs, as described in 2.3.

For the first part of the experiments, we composed our dataset from KITTI 2012 [14] and KITTI 2015 [21, 22] to simulate both calibrated and synchronized as well as the more general case of uncalibrated and unsynchronized camera array use cases. For the calibrated and synchronized camera array use case, we constructed our dataset from KITTI stereo datasets (i.e., a pair of images taken simultaneously by different cameras), consisting of unique 1578 stereo image pairs that are captured by a single pair of stereo cameras. We term this dataset as *KITTI Stereo*. By augmenting this dataset by swapping the images in the pair, hence getting a total of $1578 \times 2 = 3156$ pairs, we trained every model on 1576 image pairs, and we validated and tested every model on two different sets each with 790 image pairs from the augmented dataset.

For the second part of the experiments, we used the *Cityscape* dataset [9], consisting of 5000 stereo image pairs, where 2975 image pairs were used for training, and 500 and

1525 image pairs were used as validation and test dataset, respectively. Similarly to *KITTI Stereo*, this dataset aims to illustrate calibrated and synchronized camera array use case.

For the third part of the experiments, we simulated the general case of uncalibrated and unsynchronized camera arrays. We built the dataset from 21 stereo pairs per scene obtained sequentially from each of the 789 scenes. We name this dataset as *KITTI General*. We constructed this dataset from pairs of images, where one image is taken from the left camera and the second image from the right camera, but now, the images are taken from different time steps (unsynchronized), in our case, 1 to 3 time steps apart. Also, the images are taken up to approximately 9 meters apart (uncalibrated). This results in objects differing in scale and position between the two images, or even sometimes not appearing in one of the images at all. For this dataset, we trained, validated, and tested the models on 174936, 912, and 3607 image pairs, respectively. We evaluated the image quality performance of the models using multi-scale structural similarity index measure (MS-SSIM), which is widely reported to be a more realistic measure for human perception of image quality [35], in comparison with mean-squared error distortion. Refer to supplementary material to see sample image pairs from all datasets.

3.2. Training

For both *KITTI Stereo* and *KITTI General* datasets, we center-crop each 375×1242 image to obtain images of size 370×740 , and consequently, downsample them to 128×256 . For the *Cityscape* dataset, we directly downsample images to 128×256 . We train the benchmark models, as well as the proposed approach, with different values of λ to obtain points in different regions of the rate-distortion curves, using MS-SSIM metric for the reconstruction loss. We train all models for 500K iterations, using randomly initialized network weights. We train the models using AMS-Grad optimizer [28], with a learning rate of $1 \cdot 10^{-4}$, where we reduce the learning rate by a factor of 10 when the loss function stagnates down to a learning rate of $1 \cdot 10^{-7}$. Similarly to [24], we opt for a batch size of 1 considering the relatively small sizes of datasets under consideration. For comparison, we also train the models proposed in [24] and [2], which will be referred to as NDIC and DSIN, respectively, by using the provided codes²³. For NDIC, we used the “Ballé2017” backbone, and the model hyperparameters were kept the same. For KITTI Stereo dataset, we

²<https://github.com/ipc-lab/NDIC>,
<https://github.com/ayziksha/DSIN>.

³We did not conduct experiments with [36] since the source code of the revised version of this work is not publicly available. Furthermore, the authors mention that the exact number of channels they employ within their autoencoder network varies for different rate-distortion points, which is not provided in [36].

¹Our code is available at <https://github.com/ipc-lab/NDIC-CAM>.

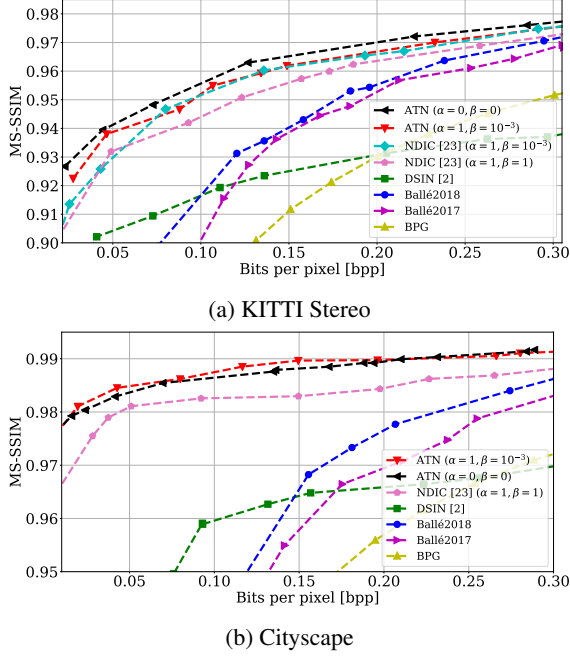


Figure 4: Comparison of different models in terms of MS-SSIM metric. “Ballé2017” and “Ballé2018” models refer to [3] and [5], respectively. “ATN” refers to our proposed approach.

used parameters ($\alpha = 1, \beta = 10^{-3}$) for the loss function and parameters ($\alpha = 1, \beta = 1$) for the rest of the experimental setup. This is due to the finding that although the parameters ($\alpha = 1, \beta = 10^{-3}$) is shown to be the best performing one considering KITTI Stereo dataset in the ablation study provided in [24], we observe that this combination of parameters induces further instability during the training process. We suspect that this is because of the reduced weighting of the regularization term controlling the complexity of the common information to be extracted (see Eq. (1)).

3.3. Experimental results

In this section, we evaluate the performance of the proposed model, which we refer to as “ATN”, and compare it with the NDIC model [24] and the DSIN model [2] (see Fig. 4). In addition to DSIN and NDIC models (discussed in Section 1.1.4), we also assess BPG as well as the DNN-aided compression schemes introduced in [3] and [5], which will be referred to as “Ballé2017” and “Ballé2018”, respectively. Following [29], we opt for 4:4:4 chroma format for BPG. It is important to remark that the point-to-point schemes such as BPG and data-driven ones such as [3, 5] do not exploit the side information at the decoder side. Looking at Fig. 4, we observe a significant improvement in performance with our proposed approach compared to the NDIC model on *KITTI Stereo* and *Cityscape* datasets. Note that in

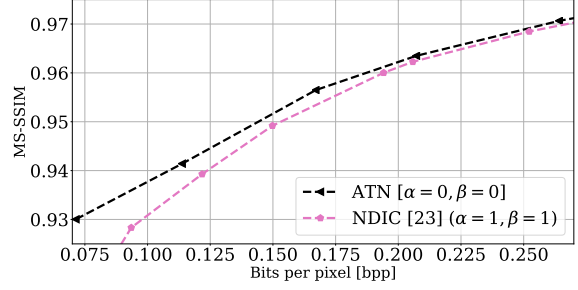


Figure 5: Comparison of the proposed approach and NDIC considering MS-SSIM metric on KITTI General dataset.

general, ATN with hyperparameters ($\alpha = 0, \beta = 0$) achieves better performance than the one with ($\alpha = 1, \beta = 10^{-3}$) on the KITTI Stereo dataset, and comparable performance to ATN with ($\alpha = 1, \beta = 10^{-3}$) on the Cityscape dataset. We argue that in order for CAMs to do feature alignment, the inputs $\hat{\mathbf{v}}_x^{(i)}$ and $\mathbf{v}_y^{(i)}$ to the CAM must be correlated. Note that this correlation is provided by the variable \mathbf{w} . By applying more pronounced regularization on \mathbf{w} , the amount of common information \mathbf{w} extracted is reduced, and $\hat{\mathbf{v}}_x^{(i)}$ and $\mathbf{v}_y^{(i)}$ become less correlated, thus reducing the efficiency of the CAMs. We also note that the proposed solution significantly improves the performance compared to DSIN in experiments with both datasets, suggesting that the proposed differentiable way of aligning the corresponding patches in the two images is better than the “search-based” patch-matching algorithm adopted by the side information (SI) finder block in [2]. We also report the results on *KITTI General* dataset in Fig. 5. The gains achieved by distributed compression models on KITTI General are notably less in comparison to those achieved on KITTI Stereo, since there is less correlation to exploit between images from different time steps for this dataset. Even in this more general setup where images are only loosely co-located in space or time, our method outperforms NDIC, where gains are more prominent in low bit rate regime.

We also provide a visual comparison of reconstructions by NDIC and our model in Fig. 6 and 7. Observe that our proposed approach captures the fine details better than NDIC, while scoring lower bit rates. Our model is especially successful in capturing the texture and color details thanks to CAM components that make use of the side information image in a superior way by aligning and morphing the corresponding patches within intermediate latents. Knowing that the objects closer to the cameras experience a larger shift from one stereo image to the other one, we can observe that the improvement in visual quality due to patch alignment done using CAMs is most evident in objects and features closer to the stereo cameras.

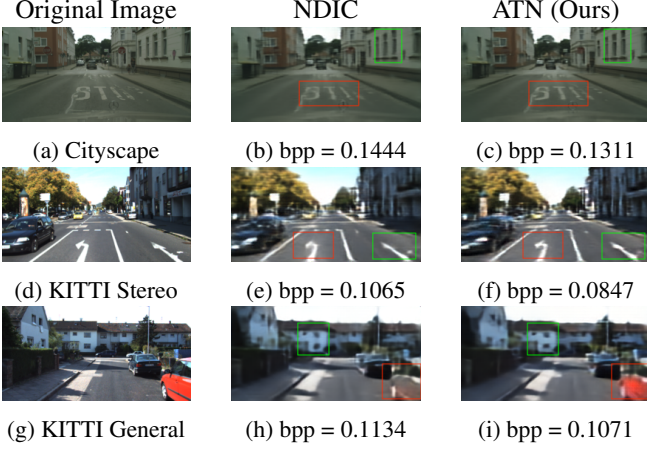


Figure 6: Visual comparison of different models trained for the MS-SSIM metric. “NDIC” refers to the model proposed in [24].

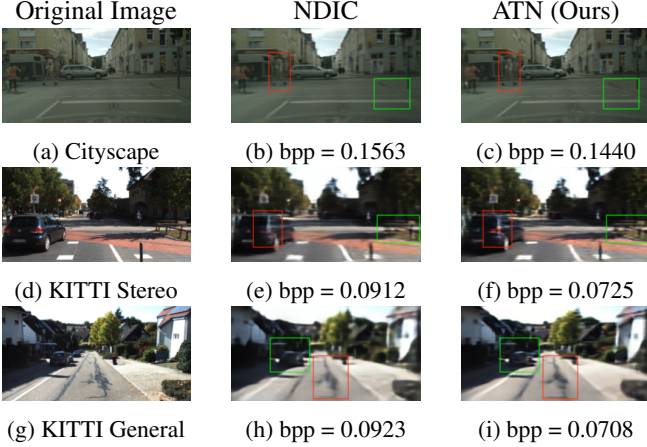


Figure 7: Additional examples for visual comparison of different models trained for the MS-SSIM metric.

3.3.1 Feature alignment

In Fig. 8, a sample channel from the latent feature representation $\hat{\mathbf{v}}_x^{(2)}$ after the second layer of the decoder \mathbf{g}_{sx} , and the corresponding channel from the output of the CAM, are shown. Observe that the road edges in the bottom left corner of the original left and right images are at different locations, but after the application of CAM to $\hat{\mathbf{v}}_x^{(2)}$ and $\mathbf{v}_y^{(2)}$, the features corresponding to the road edge in the CAM output $\mathbf{v}_{CAM}^{(2)}$ are aligned with those in $\hat{\mathbf{v}}_x^{(2)}$. This indicates that the CAM layer learns how to align the features in the latent representation of the SI with those in the latent representation of the input image, allowing more efficient utilization of the features available in the SI.

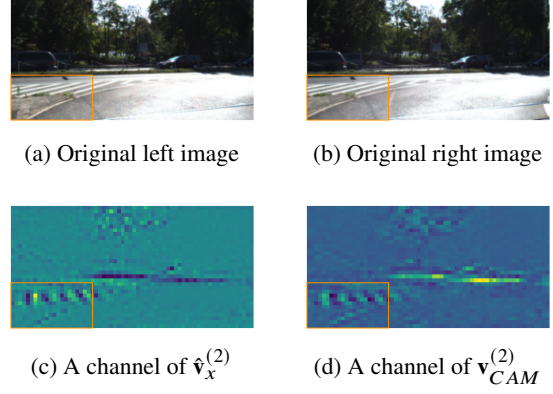


Figure 8: Alignment of feature maps.

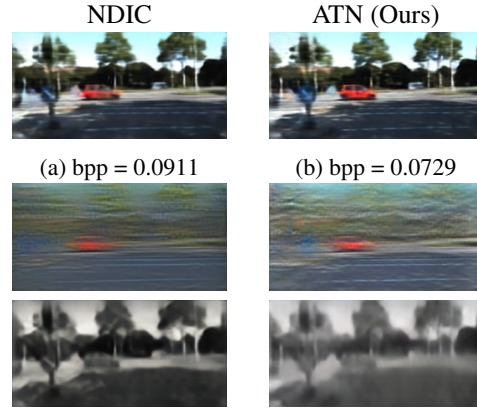


Figure 9: For a similar reconstruction quality in KITTI Stereo dataset (1st row), the decomposition of common information (2nd row) and private information (3rd row) for NDIC and for our proposed approach.

3.3.2 Visualization of private and common information

In Fig. 9, we provide visualizations of the private and common information components obtained for NDIC as well as our model. We generate the private information visualization by plotting the output of the decoder when the side information image is replaced with a fixed array of 0.5. This is done in order to block any relevant information that the decoder might extract from the SI. We also generate the common information visualization by plotting the output of the decoder when the input image is replaced with a fixed array of 0.5, in order to block all information from the input image. Consistent with [24], we observe that the common information mostly captures the global color and texture details whereas private information captures the structural content (e.g., objects and edges). For a similar reconstruction quality, observe that our approach yields a richer and more defined common information, and yields lower fi-

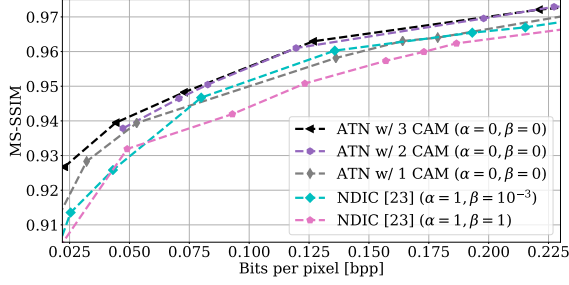


Figure 10: Ablation study experiments on the number of “CAM” layers on the proposed architecture, using the MS-SSIM metric on the KITTI Stereo dataset.

delity private information compared to NDIC. This explains why our model is able to capture finer details compared to NDIC, while scoring lower bit rates, which depends on the fidelity of the private information sent by the encoder. By extracting more common information from the side information image at the decoder side, the proposed approach relies less on the information transmitted from the encoder to achieve a similar reconstruction quality.

3.3.3 Ablation study

The outputs of each layer of the decoders capture features at different scales, where the initial layers capture large scale features, and the later layers capture the small scale features. Therefore, CAMs applied to the outputs of the initial layers do large scale alignment, while CAMs applied to the later layers do alignment of the small-scale features. To study the impact of each CAM component in our approach (see Fig. 2 for the baseline architecture), we carry out an ablation study on the number of CAMs, and compare the performances in Fig. 10. We remove the CAM layers starting from the last convolutional layer, moving in the direction of the first layer. As seen in the plot, removing 1 CAM layer does not affect the performance significantly. However, removing the second CAM layer results in a significant reduction in performance. See Fig. 11 for a visual comparison of the performances between the model with 1 CAM and the model with 3 CAM components.

3.3.4 Drawbacks

We discuss a few drawbacks and limitations of this work. Like most deep learning-based image compression works, our method is dataset-dependent, that is, it performs well on the data distribution it is trained on, while not guaranteeing a good performance on images from another distribution. Another limitation is that the proposed model has almost double the number of parameters compared to [24], thus resulting in slower inference times than other approaches.



Figure 11: Reconstructed images obtained for the model having only 1 CAM and 3 CAM components. Having more CAM layers helps the model to preserve finer details while scoring a lower bit rate.

Please refer to the supplementary materials for a discussion of number of model parameters and inference times.

4. Conclusion

We presented a new method for distributed stereo image compression, which makes use of cross-attention mechanisms in order to align the feature maps of the intermediate layers in the decoding stage. The method achieves a superior performance in exploiting the correlation between the decoder-only side information image and the image to be reconstructed, compared to the solution provided in [24]. We have shown that this approach achieves good reconstruction quality even at very low bit regimes, substantially outperforming the single image compression models, as well as surpassing the previous works on distributed image compression with side information. Even for a more general camera array use case with uncalibrated and unsynchronized images, we have shown that the proposed method is on par or superior in performance with respect to the approach in [24]. The ablation study shows that there is diminishing marginal benefit with the increasing number of CAM components employed in the decoding pipeline, which provides a trade-off between decoding complexity and performance.

5. Acknowledgements

This work received funding from the European Research Council (ERC) through Starting Grant BEACON (no. 677854) and from the UK EPSRC (project CONNECT with grant no. EP/T023600/1 and project SONATA with grant no. EP/W035960/1).

References

- [1] E. Agustsson, M. Tschannen, F. Mentzer, R. Timofte, and L. V. Gool. Generative adversarial networks for extreme learned image compression. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 221–231, 2019.
- [2] S. Ayzik and S. Avidan. Deep image compression using decoder side information. In *European Conference on Computer Vision (ECCV)*, 2020.
- [3] J. Ballé, V. Laparra, and E. P. Simoncelli. End-to-end optimized image compression. In *International Conference on Learning Representations (ICLR)*, 2017.
- [4] J. Ballé, V. Laparra, and E. P. Simoncelli. Density modeling of images using a generalized normalization transformation. In *International Conference on Learning Representations (ICLR)*, 2016.
- [5] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston. Variational image compression with a scale hyperprior. In *International Conference on Learning Representations (ICLR)*, 2018.
- [6] F. Bellard. BPG image format. <https://bellard.org/bpg/>, 2014.
- [7] T. Chen, H. Liu, Zhan Ma, Q. Shen, X. Cao, and Y. Wang. End-to-end learnt image compression via non-local attention optimization and improved context modeling. *IEEE Transactions on Image Processing*, 30:3179–3191, 2021.
- [8] Z. Cheng, H. Sun, M. Takeuchi, and J. Katto. Learned image compression with discretized gaussian mixture likelihoods and attention modules. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7936–7945, 2020.
- [9] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [10] X. Deng, W. Yang, R. Yang, M. Xu, E. Liu, Q. Feng, and R. Timofte. Deep homography for efficient stereo image compression. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1492–1501, 2021.
- [11] E. Diao, J. Ding, and V. Tarokh. Drasic: Distributed recurrent autoencoder for scalable image compression. In *Data Compression Conference (DCC)*, 2020.
- [12] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- [13] Peter Gacs and J. Körner. Common information is far less than mutual information. *Problems of Control and Information Theory*, 2, 01 1973.
- [14] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [15] Z. Guo, Z. Zhang, R. Feng, and Z. Chen. Causal contextual prediction for learned image compression. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(4):2329–2341, 2022.
- [16] J. Lee, S. Cho, and S.K. Beack. Context-adaptive entropy model for end-to-end optimized image compression. In *International Conference on Learning Representations (ICLR)*, 2019.
- [17] J. Lei, X. Liu, B. Peng, D. Jin, W. Li, and J. Gu. Deep stereo image compression via bi-directional coding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19669–19678, June 2022.
- [18] M. Li, W. Zuo, S. Gu, J. You, and D. Zhang. Learning content-weighted deep image compression. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(10):3446–3461, 2021.
- [19] J. Liu, S. Wang, and R. Urtasun. Dsic: Deep stereo image compression. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3136–3145, 2019.
- [20] F. Mentzer, G. Toderici, M. Tschannen, and E. Agustsson. High-fidelity generative image compression. In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS’20*, Red Hook, NY, USA, 2020. Curran Associates Inc.
- [21] M. Menze, C. Heipke, and A. Geiger. Joint 3d estimation of vehicles and scene flow. In *ISPRS Workshop on Image Sequence Analysis (ISA)*, 2015.
- [22] M. Menze, C. Heipke, and A. Geiger. Object scene flow. *ISPRS Journal of Photogrammetry and Remote Sensing (JPRS)*, 2018.
- [23] D. Minnen, J. Ballé, and G. D. Toderici. Joint autoregressive and hierarchical priors for learned image compression. In *Advances in Neural Information Processing Systems*, 2018.
- [24] N. Mital, E. Özyılkan, A. Garjani, and D. Gündüz. Neural distributed image compression using common information. In *2022 Data Compression Conference (DCC)*, pages 182–191, 2022.
- [25] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [26] Y. Patel, S. Appalaraju, and R. Manmatha. Saliency driven perceptual image compression. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 227 – 236, 2021.
- [27] P. Ramachandran, N. Parmar, A. Vaswani, I. Bello, A. Levskaya, and J. Shlens. *Stand-Alone Self-Attention in Vision Models*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- [28] S. J. Reddi, S. Kale, and S. Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations (ICLR)*, 2018.

- [29] O. Rippel and L. Bourdev. Real-time adaptive image compression. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.
- [30] A. Skodras, C. Christopoulos, and T. Ebrahimi. The jpeg 2000 still image compression standard. *IEEE Signal Processing Magazine*, 2001.
- [31] D. Slepian and J. Wolf. Noiseless coding of correlated information sources. *IEEE Transactions on Information Theory*, 1973.
- [32] L. Theis, W. Shi, A. Cunningham, and F. Huszár. Lossy image compression with compressive autoencoders. *ArXiv*, abs/1703.00395, 2017.
- [33] G. Toderici, D. Vincent, N. Johnston, S. J. Hwang, D. Minnen, J. Shor, and M. Covell. Full resolution image compression with recurrent neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [34] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need, 2017.
- [35] Z. Wang, E. Simoncelli, and A. Bovik. Multiscale structural similarity for image quality assessment. In *Asilomar Conference on Signals, Systems and Computers*, 2003.
- [36] J. Whang, A. Acharya, H. Kim, and A. G. Dimakis. Neural distributed source coding. <https://arxiv.org/abs/2106.02797>, 2021.
- [37] H. S. Witsenhausen. On sequences of pairs of dependent random variables. *Siam Journal on Applied Mathematics*, 28:100–113, 1975.
- [38] M. Wödlinger, J. Kotera, J. Xu, and R. Sablatnig. Sasic: Stereo image compression with latent shifts and stereo attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 661–670, June 2022.
- [39] A. Wyner and J. Ziv. The rate-distortion function for source coding with side information at the decoder. *IEEE Transactions on Information Theory*, 1976.
- [40] X. Ying, Y. Wang, L. Wang, W. Sheng, W. An, and Y. Guo. A stereo attention module for stereo image super-resolution. *IEEE Signal Processing Letters*, 27:496–500, 2020.