
Task-aware Distributed Source Coding under Dynamic Bandwidth

Po-han Li^{*1} Sravan Kumar Ankireddy^{*1} Ruihan Zhao¹ Hossein Nourkhiz Mahjoub² Ehsan Moradi-Pari²
Ufuk Topcu¹ Sandeep Chinchali¹ Hyeji Kim¹

Abstract

Efficient compression of correlated data is vital for reducing communication overload in multi-sensor networks, where sensors independently compress data for transmission to a central node. To optimize performance, it's crucial for the compressor to learn only task-relevant features, considering the fluctuating bandwidth availability. Our work introduces a novel distributed compression framework, Neural Distributed Principal Component Analysis (NDPCA), comprising independent encoders and a joint decoder. NDPCA adapts flexibly to varying bandwidth, reducing computational and storage demands by employing a single model. By learning low-rank task representations and efficiently allocating bandwidth among sensors, NDPCA achieves a balanced trade-off between performance and bandwidth utilization. Experimental results demonstrate NDPCA's effectiveness, improving success rates in multi-view robotic arm manipulation by 9% and enhancing object detection accuracy in satellite imagery tasks by 14% compared to an autoencoder with uniform bandwidth allocation.¹

1. Introduction

Efficient data compression is pivotal in multi-sensor networks to minimize communication overload. Due to the limited bandwidth of such networks, it is often impractical to transmit all sensor data to a central server, and compressing data is necessary. In many cases, the sensors, so-called sources, observe correlated data, which are only processed by a downstream task, *e.g.*, an object detection model, but not by human eyes. For example, satellites observe overlap-

^{*}Equal contribution ¹The University of Texas at Austin, USA ²Honda Research Institute, USA. Correspondence to: Po-han Li <pohanli@utexas.edu>, Sravan Kumar Ankireddy <sra-
van.ankireddy@utexas.edu>.

Presented at 37th Conference on Neural Information Processing Systems (NeurIPS 2023)

¹<https://github.com/UTAustin-SwarmLab/Task-aware-Distributed-Source-Coding>.

ping images and transmit them through limited bandwidth to a central server on Earth. Hence, sources should not transmit redundant information from correlated data and only transmit features relevant to the downstream task. It is important to compress each source independently to reduce the communication overload in the network. Literature refers to this setting as distributed source coding. Together, we name the distributed compression of task-relevant features *task-aware distributed source coding*.

However, existing compression methods fail to combine three aspects: 1. Existing distributed compression methods perform poorly in the presence of a task model. Although neural networks have been shown to be capable of compressing stereo images (Ballé et al., 2016; 2018) and correlated images (Zhang et al., 2023), existing methods focus on reconstructing data, but not for downstream tasks. 2. Existing task-aware compression methods cannot take advantage of the correlation of sources. Previous works only consider compressing task-relevant features of single source (Cheng et al., 2022a; Ji et al., 2012; Choi and Han, 2020; Nakanoya et al., 2023; Cheng et al., 2021), but not multiple correlated sources. 3. All existing methods, especially those based on neural networks, only compress data to a fixed compression level but not multiple levels. Thus, they cannot operate in environments with different demands of compression levels. We use the term bandwidth to indicate the information bottleneck in the dimension of transmitted data, and more related works are discussed in Appendix A.

We present neural distributed principal component analysis (NDPCA), a distributed compression framework that transmits task-relevant features at multiple compression levels. Fig. 1 illustrates the scenario where the central node requires data from all sources, and network bandwidth varies over time. NDPCA consists of neural encoders E_1, E_2, \dots, E_K that independently compress correlated data X_1, X_2, \dots, X_K to latent representations Z_1, Z_2, \dots, Z_K . The distributed principal component analysis (DPCA) module compresses these representations to any dimension based on the current bandwidth. At the central node, a neural decoder reconstructs the representations $\hat{Z}_1, \hat{Z}_2, \dots, \hat{Z}_k$ to $\hat{X}_1, \hat{X}_2, \dots, \hat{X}_K$ and feeds them into a task. NDPCA combines a neural autoencoder and the DPCA module to generate task-relevant compressible

representations.

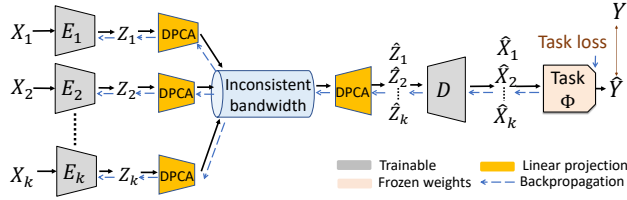


Figure 1: **Task-aware distributed source coding with NDP-CA.** X_1, \dots, X_k are correlated data sources. Encoders E_1, \dots, E_k independently compress data to latent representations Z_1, \dots, Z_k . The DPCA module allocates the bandwidth of sources based on the importance of the task Φ . The goal is to find the optimal encoders and corresponding decoder that minimize the final task loss.

Contributions: First, we formulate the task-aware distributed source coding problem (Sec. 2). Second, we provide a theoretical justification for the framework by analyzing the case of a linear compressor and a task (Sec. 3). Third, we propose a task-aware distributed source coding framework, NDP-CA, that learns a single model for different compression levels (Sec. 4). We validate NDP-CA with an object detection task of satellite imagery (Sec. 5), resulting in a 14% increase in accuracy compared to an autoencoder with uniform bandwidth allocation.

2. Problem Formulation

Consider a set of K correlated sources. Let $x_i \in \mathbb{R}^{n_i}$ denote the sample from source i where $i \in \{1, 2, \dots, K\}$. Samples from each source are compressed independently by encoder E_i to a latent representation $z_i \in \mathbb{R}^{m_i}$ such that $\sum_{i=1}^K m_i = m$, where m is the total bandwidth available. A joint decoder D receives the representations $\{z_1, z_2, \dots, z_k\}$ and reconstructs the sources $\{\hat{x}_1, \hat{x}_2, \dots, \hat{x}_k\} = \{D(E_1(x_1)), D(E_2(x_2)), \dots, D(E_k(x_k))\}$. In the presence of a task Φ , it takes the reconstructed inputs to compute the final output $\Phi(\hat{x}_1, \hat{x}_2, \dots, \hat{x}_k)$. The goal is to find a set of encoders and a decoder such that the task loss $\mathcal{L}_{\text{task}}$ is minimized. We call this problem as *task-aware distributed source coding*, which is the focus of this paper:

$$\begin{aligned} & \underset{E_1, \dots, E_k, D}{\operatorname{argmin}} \quad \mathcal{L}_{\text{task}}(Y, \hat{Y}) \\ & \text{s.t.} \quad Y = \Phi(x_1, \dots, x_k), \hat{Y} = \Phi(\hat{x}_1, \dots, \hat{x}_k) \\ & \quad \quad \quad (\text{Task-aware distributed source coding}), \end{aligned} \quad (1)$$

where $\mathcal{L}_{\text{task}}$ is the task loss, *e.g.*, the difference of bounding boxes when the task is object detection.

Bandwidth allocation: In the previous formulations, we assume that the output dimensions of encoders are known a priori. However, the dimensions determine the compression of each encoder, which is also a design factor. Given the

total available bandwidth m , we first need to obtain the optimal m_i for each source i ; then, we can design the optimal encoders and decoder accordingly. Finding the optimal set of bandwidths for a given task is a long-standing open problem, even for the simple task of a modulo-two sum of two binary sources (Korner and Marton, 1979). Also, existing works (Zhang et al., 2023; Whang et al., 2021; Mital et al., 2023) largely assume a fixed latent dimension for sources and train different models for different total available bandwidth m , which is, of course, suboptimal. NDP-CA provides heuristics to the underlying key challenge of optimally allocating available bandwidth, *i.e.*, deciding m_i , while adapting to different total bandwidths m with a single model.

3. Theoretical Analysis

We start with a motivating example of task-aware distributed source coding under the constraint of linear encoders, a decoder, and a linear task.

DPCA: We consider a linear task for two sources, defined by the task matrix $\Phi \in \mathbb{R}^{p \times (n_1 + n_2)}$, where the sources $x_1 \in \mathbb{R}^{n_1}$ and $x_2 \in \mathbb{R}^{n_2}$ are of dimensions n_1 and n_2 , respectively, and the task output is given by $y = \Phi x \in \mathbb{R}^p$, where $x = [x_1^\top, x_2^\top]^\top$. Without loss of generality, we assume the sources to be zero-mean. Now, we have N observations of two sources $X_1 \in \mathbb{R}^{n_1 \times N}$ and $X_2 \in \mathbb{R}^{n_2 \times N}$ and their corresponding task outputs $Y = \Phi(X) \in \mathbb{R}^{p \times N}$, where $X = [X_1^\top, X_2^\top]^\top$. We aim to design the optimal linear encoding matrices (encoders) $E_1 \in \mathbb{R}^{m_1 \times n_1}$, $E_2 \in \mathbb{R}^{m_2 \times n_2}$, and the decoding matrix (decoder) $D \in \mathbb{R}^{(n_1 + n_2) \times (m_1 + m_2)}$ that minimizes the task loss defined as the Frobenius norm of $\Phi(X) - \Phi(\hat{X})$, where \hat{X} is the reconstructed X . For now, we assume that m_1 and m_2 are given. Letting $Z_1 = E_1 X_1 \in \mathbb{R}^{m_1 \times N}$ and $Z_2 = E_2 X_2 \in \mathbb{R}^{m_2 \times N}$ denote the encoded representations and $M = \Phi D$ denote the product of the task and decoder matrices, we solve the optimization problem:

$$E_1^*, E_2^*, M^* = \underset{E_1, E_2, M}{\operatorname{argmin}} \quad \|Y - MZ\|_2^2 \quad (2a)$$

$$\text{s.t.} \quad Z = \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} = \begin{bmatrix} E_1 X_1 \\ E_2 X_2 \end{bmatrix}, \quad (2b)$$

$$ZZ^\top = \mathbb{I}_m, \quad \hat{Y} = MZ, \quad Y = \Phi \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \quad (2c)$$

Note that solving M is identical to solving the decoder D since we can always convert M to D by the generalized inverse of task Φ . We constrain the representations to be orthonormal vectors in (2c) as in the normalization in principal component analysis (PCA) for the compression of a single source (Jolliffe, 2011).

We discuss the detailed solution of DPCA in Appendix B

and provide a brief summary here. First, a preprocessing step removes the correlation part of X_1 from X_2 by subtracting the least-square estimator $\hat{X}_2(X_1)$:

$$\tilde{X}_2 = X_2 - \hat{X}_2(X_1) = X_2 - X_2 X_1^\top (X_1 X_1^\top)^{-1} X_1. \quad (3)$$

The orthogonality principle of least-square estimators (Kay, 1993) ensures that $X_1 \tilde{X}_2^\top = \mathbf{0}_{n_1 \times n_2}$. Thus, we can decouple the original problem in (2) into subproblems where each corresponds to a single encoder, which can be solved by canonical correlation analysis (Hardoon et al., 2004).

Dynamic bandwidth: We extend our approach to determine the optimal bandwidth allocation given a total bandwidth m . By solving DPCA with $m_1 = n_1$ and $m_2 = n_2$, we obtain optimal encoders E_1^* and E_2^* , as well as pairs of canonical directions and correlations. Similar to PCA, these pairs can be seen as a generalization of singular vectors and values, and the sums of squares of canonical correlations are the optimal values of (2). We select the largest m pairs of canonical correlations and directions. These canonical correlations determine the optimal encoders E_1^* , E_2^* , and decoder D^* . The importance of a direction to the task is indicated by the canonical correlations, so we prioritize the transmission of important directions.

Performance analysis of DPCA: When DPCA compresses new data matrices with encoder E_1^* and E_2^* , the preprocessing step (3) is invalid as the encoders cannot communicate with each other. So for DPCA to perform optimally while skipping the step, the two data matrices need to be uncorrelated, namely, $\hat{X}_2(X_1) = 0$, because in such case, the preprocessing step removes nothing from the data sources. Given that correlated sources lead to suboptimality of DPCA, we characterize the performance under the same bandwidth in Lemma C.1 with the simplest case of reconstruction, namely, $\Phi = \mathbb{I}_p$. In this case, the canonical correlation analysis is relaxed to the singular value decomposition, which is later used for NDPCA in Sec. 4. Lemma C.1 concludes that DPCA performs more closely the optimal joint compression, PCA, as the covariance decreases.

4. Neural Distributed Principal Component Analysis

Theoretical analysis reveals that DPCA has limitations: it optimally compresses data only when sources are uncorrelated and is limited to linear tasks but can allocate the bandwidth dynamically. Neural autoencoders, on the other hand, excel at fixed-dimension compression but lack dynamic bandwidth allocation. To address this contrast, we propose neural distributed principal component analysis (NDPCA), which combines a neural autoencoder and DPCA. This integrated approach enables efficient compression and optimal bandwidth allocation. NDPCA has two encoding stages, as shown in Fig. 1: First, the neural encoder at each i -th source

encodes data X_i to a fixed-dimensional representation Z_i for $i \in [K]$. Then, DPCA adapts the dimension of Z_i via linear matrices according to the available bandwidth as per Sec. 3. Similarly, the decoding of NDPCA is also performed in two stages. First, the DPCA linear decoder reconstructs the K fixed-dimensional representations $\hat{Z}_1, \dots, \hat{Z}_K$, based on which the joint neural decoder generates the estimate of data $\hat{X}_1, \dots, \hat{X}_K$. These estimates are then passed to the neural task model Φ to obtain the final task output \hat{Y} . Since we have a non-linear task model here, DPCA mainly adapts the dimension appropriately as needed; the role of the DPCA here is to reliably reconstruct the embedding \hat{Z} , which corresponds to the case described in Lemma C.1 with the task matrix Φ as identity.

Training procedure: During the training of NDPCA, we assume the task model is pre-trained, and we do not update its weights. We aim to learn the K neural encoders and the joint neural decoder, which minimize the loss function:

$$\mathcal{L}_{\text{tot}} = \lambda_{\text{task}} \underbrace{\|\hat{Y} - Y\|_F^2}_{\text{task loss}} + \lambda_{\text{rec}} \underbrace{\sum_{i=1}^K \|\hat{X}_i - X_i\|_F^2}_{\text{reconstruction loss}}. \quad (4)$$

In the task-aware setting when $\lambda_{\text{rec}} = 0$, the neural autoencoder fully restores task-relevant features, which is the main focus of this paper. When $\lambda_{\text{task}} = 0$, the neural autoencoder learns to reconstruct the data X , which is the task-agnostic setting later compared in Sec. 5.

To encourage NDPCA to work well under *various available bandwidths* with DPCA during the training phase, we need uncorrelatedness from the limitations of the DPCA. To compress representations with a few singular vectors and make NDPCA more bandwidth efficient, we need *linear compressibility*. That is, encouraging the neural autoencoder to generate low-rank representations. We tried to explicitly encourage the desired properties with additional terms in (4), but they all adversely affect the task performance. We tried to use the cosine similarity to generate uncorrelated representations as per (Bardes et al., 2022; Singh et al., 2022; Bousmalis et al., 2016; Chen et al., 2018), and the convex low-rank approximation, nuclear norm, to increase linear compressibility, as per (Salzmann et al., 2010; Fazel, 2002). For the comparison of the resulting performance, see Appendix H.1. In this regard, we propose a novel linear compression module that allows us to adapt to DPCA during training rather than using additional terms in the loss. We introduce a *random-dimension* DPCA projection module to improve performance in lower bandwidths. It projects representations Z to a random low dimension, simulating projections in various available bandwidths during inference. It can be interpreted as a differentiable singular value decomposition with a random dimension, described in Alg. 1. Note that no retraining is needed for different bandwidths,

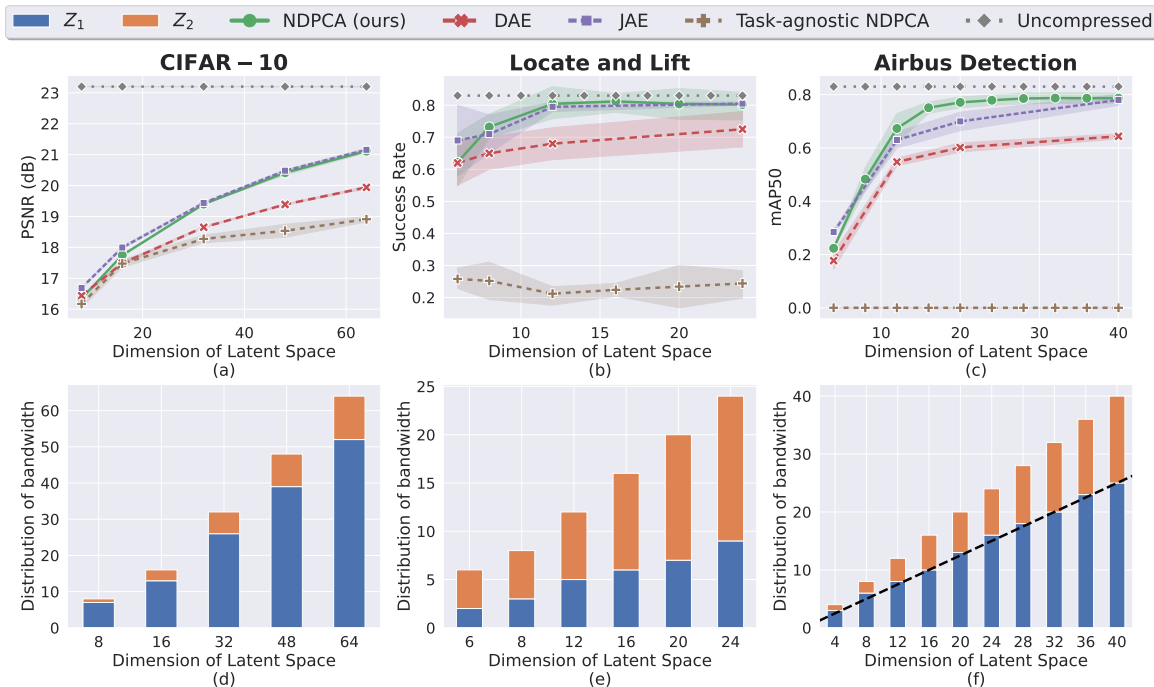


Figure 2: **Top:** Performance Comparison for 3 different tasks. Our method achieves equal or higher performance than other methods. **Bottom:** Distribution of total available bandwidth (latent space) among the two views for NDPCA (ours). The unequal allocation highlights the difference in the importance of the views for a given task.

and only the storage of a neural autoencoder and a linear matrix at each encoder and decoder is needed.

5. Experiments

We consider three different tasks to test our framework: (a) the denoising of CIFAR-10 images (Krizhevsky, 2009), (b) multi-view robotic arm manipulation (Zhan et al., 2022), which we refer to as the *locate and lift* task, and (c) object detection on satellite imagery (Defense and Intelligence, 2021). Across all the experiments, we assume that there are two data sources, referred to as *views*, each containing partial information relevant to the task. We present our results based on the testing set and refer to our proposed method, task-aware NDPCA, as NDPCA for simplicity. NDPCA includes a single autoencoder with a large dimension of representations $Z \in \mathbb{R}^{2 \times m_{\max}}$. It then compresses representations and allocates bandwidth via DPCA, as discussed previously. We show that NDPCA can bridge the performance gap between distributed autoencoders and joint autoencoders, defined below, to allocate bandwidth and avoid transmitting task-irrelevant features.

Baselines: We compare NDPCA against three baselines: (a) Task-aware joint autoencoder (JAE), where a single pair of encoder and decoder compresses both views. JAE is considered an upper bound of NDPCA since it can leverage the correlation between both views while avoiding encoding

redundant information. (b) Task-aware vanilla distributed autoencoder (DAE), where two encoders independently encode one view to equal bandwidths, and a joint decoder decodes the data. DAE is considered a lower bound of NDPCA since both encoders utilize the same bandwidth regardless of the importance of the views for the task, while NDPCA allocates bandwidths in a task-aware manner. (c) Task-agnostic NDPCA differs from NDPCA in the training loss of reconstructing the original views.

CIFAR-10 Denoising: We tackle CIFAR-10 image denoising using two noisy observations of the same image (Fig. 3 (a)). This serves as a simplified example to demonstrate NDPCA’s advantage when source importance varies. With the task requiring only 4 bits for classification, denoising is chosen to showcase performance across different bandwidths. Each observation’s significance depends solely on its noise level. View 1 is corrupted with AWGN variance 0.1^2 , while view 2 has higher corruption with variance 1. Images are normalized to $[0, 1]$ before noise addition. After compressing the noisy data, we pass the reconstructed images through a pre-trained denoising network, calculating PSNR relative to the clean image. Unequal noise levels across views result in varying task importance, suggesting a non-uniform bandwidth allocation strategy and emphasizing NDPCA’s advantage. Optimal bandwidth distribution is crucial to address the CEO problem (Berger et al., 1996; Prabhakaran et al., 2004), ensuring all views are leveraged,

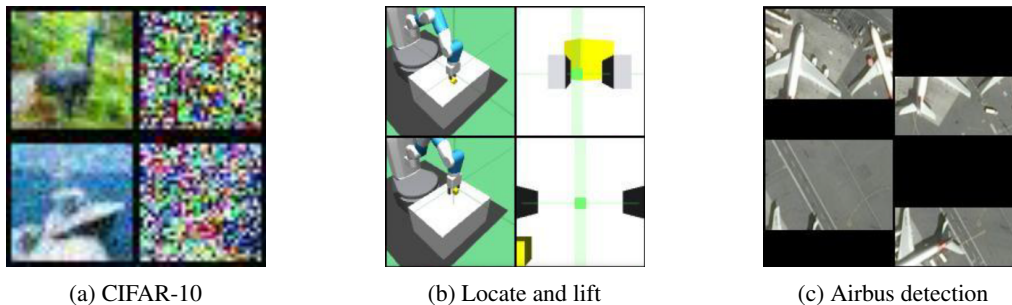


Figure 3: **Datasets:** (column 1) view 1. (column 2) view 2. In all experiments, both views are correlated, but one view is more important than the other as it contains more information relevant to the task.

even if heavily corrupted, rather than allocating zero bandwidth.

Locate and lift: For the manipulation task, we consider a scenario in which a simulated 6 degrees-of-freedom robotic arm controlled by a reinforcement learning agent inputs two camera views to locate and lift a yellow brick. We call the view from the robotic arm "arm-view" and the one recording the whole desk "side-view", as shown in Fig. 3 (b). The two views are complementary to completing the task, details discussed in Appendix H.3. We trained the agent in a supervised-learning manner. We collected a dataset of observation and action pairs (Zhao et al., 2023) and trained an agent from the dataset. Then, we defined task loss as the L_2 norm of actions from images with and without compression and trained NDPCA to minimize the task loss through the agent. Literature calls this training method "behavior cloning" (Torabi et al., 2018) as it learns from demonstrations. Behavior cloning causes a drop in performance, but this paper only focuses on the performance degradation caused by compression, so we treat the behavior cloning agent with uncompressed views as the upper bound of our method.

Airbus detection: This task considers using satellite imagery to locate Airbuses. Satellites observe overlapping images of an airport and transmit data to Earth through limited bandwidth, as shown in Fig. 3 (c). We crop all images in the dataset into smaller pieces (224×224 pixels). The two data sources are the upper 160 pixels (source 1) and the lower 104 pixels of the image (source 2) with 40 pixels overlapped. Our object detection model follows the paper "You Only Look Once" (Yolo) (Redmon et al., 2016). The task loss here is the difference between object detection loss with and without compression.

Results: Our key findings are: (1) Task-aware NDPCA outperforms task-agnostic NDPCA, and (2) bandwidth allocation should correlate with task importance.

Across experiments (Fig. 2(a)-(c)), task-aware NDPCA notably surpasses task-agnostic NDPCA and DAE, which evenly distribute bandwidth. Task-aware NDPCA exhibits

graceful performance degradation with available bandwidth, requiring no additional training or storage of multiple models, unlike DAE and JAE, which necessitate retraining for each compression level.

In Fig. 2(a), denoising CIFAR-10 with NDPCA trained at $(m_{\min}, m_{\max}) = (8, 64)$ reveals NDPCA's ability to flexibly allocate bandwidth, yielding a 1.2 dB PSNR gain over DAE when $m = 64$.

Fig. 2(b) showcases locate-and-lift task results with NDPCA trained at $(m_{\min}, m_{\max}) = (8, 48)$. NDPCA achieves a 9% higher success rate compared to DAE when $m = 24$ due to its effective allocation of bandwidth, prioritizing the more crucial arm view.

In Fig. 2(c), Airbus detection with NDPCA trained at $(m_{\min}, m_{\max}) = (8, 40)$ demonstrates up to a 14% gain in mAP50 over DAE. Interestingly, NDPCA's empirical bandwidth allocation closely aligns with the theoretical ratio, suggesting its proficiency in capturing task importance.

Comparison with JAE: Despite JAE's theoretical advantage, NDPCA outperforms not only DAE but also JAE in some cases. This superiority is attributed to NDPCA's more efficient representation in higher-dimensional latent space compared to JAE's direct learning of low-dimensional representation. This approach is akin to pruning large neural networks to identify effective sparse sub-networks.

6. Conclusion

We proposed a theoretically-grounded linear distributed compressor, DPCA, and analyzed its performance compared to the optimal joint compressor. Then, we designed a distributed compression framework called NDPCA by combining a neural autoencoder and DPCA to allocate bandwidth according to their importance to the task. Experiments on Airbus detection showed that NDPCA near-optimally outperforms task-agnostic or equal-bandwidth compression schemes. Moreover, NDPCA requires only a single model for different compression levels, which makes it suitable for settings with dynamic bandwidths.

References

- Ultralytics yolov8 docs. <https://docs.ultralytics.com/>, 2023. [Online; accessed 16-April-2023].
- J. Ballé, V. Laparra, and E. P. Simoncelli. End-to-end optimized image compression. *arXiv preprint arXiv:1611.01704*, 2016.
- J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston. Variational image compression with a scale hyperprior. *arXiv preprint arXiv:1802.01436*, 2018.
- A. Bardes, J. Ponce, and Y. LeCun. VICReg: Variance-invariance-covariance regularization for self-supervised learning. In *International Conference on Learning Representations*, 2022.
- T. Berger, Z. Zhang, and H. Viswanathan. The CEO problem. *IEEE Transactions on Information Theory*, 42(3):887–902, 1996. doi: 10.1109/18.490552.
- K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, and D. Erhan. Domain separation networks. In *Advances in Neural Information Processing Systems*, 2016.
- A. Buslaev, V. I. Iglovikov, E. Khvedchenya, A. Parinov, M. Druzhinin, and A. A. Kalinin. Alumentations: Fast and flexible image augmentations. *Information*, 11(2), 2020. ISSN 2078-2489. doi: 10.3390/info11020125.
- R. T. Q. Chen, X. Li, R. B. Grosse, and D. K. Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *Advances in Neural Information Processing Systems*, 2018.
- J. Cheng, M. Pavone, S. Katti, S. Chinchali, and A. Tang. Data sharing and compression for cooperative networked control. In *Advances in Neural Information Processing Systems*, 2021.
- J. Cheng, S. Chinchali, and A. Tang. Task-aware network coding over butterfly network. *arXiv preprint arXiv:2201.11917*, 2022a.
- J. Cheng, A. Tang, and S. Chinchali. Task-aware privacy preservation for multi-dimensional data. In *Proceedings of the 39th International Conference on Machine Learning*, 2022b.
- J. Choi and B. Han. Task-aware quantization network for jpeg image compression. In *European Conference on Computer Vision*, 2020. ISBN 978-3-030-58565-5.
- A. Defense and S. Intelligence. Airbus aircraft detection. <https://www.kaggle.com/datasets/airbusgeo/airbus-aircrafts-sample-dataset>, 2021. [Online; accessed 03-April-2023].
- E. Diao, J. Ding, and V. Tarokh. DRASIC: Distributed recurrent autoencoder for scalable image compression. In *Data Compression Conference*, 2020. doi: 10.1109/DCC47342.2020.00008.
- Y. Dubois, B. Bloem-Reddy, K. Ullrich, and C. J. Maddison. Lossy compression for lossless prediction. In *Advances in Neural Information Processing Systems*, 2021.
- M. Fazel. *Matrix rank minimization with applications*. PhD thesis, Stanford University, 2002.
- J. Frankle and M. Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2019.
- P. han Li, S. P. Chinchali, and U. Topcu. Differentially private timeseries forecasts for networked control. In *American Control Conference*, 2023.
- D. R. Hardoon, S. Szedmak, and J. Shawe-Taylor. Canonical correlation analysis: An overview with application to learning methods. *Neural Computation*, 16(12):2639–2664, 2004. doi: 10.1162/0899766042321814.
- I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner. beta-VAE: Learning basic visual concepts with a constrained variational framework. In *International Conference on Learning Representations*, 2017.
- N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, and S. Gelly. Parameter-efficient transfer learning for NLP. In *International Conference on Machine Learning*, 2019.
- E. J. Hu, yelong shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- R. Ji, H. Yao, W. Liu, X. Sun, and Q. Tian. Task-dependent visual-codebook compression. *IEEE Transactions on Image Processing*, 21(4):2282–2293, 2012. doi: 10.1109/TIP.2011.2176950.
- I. Jolliffe. *Principal Component Analysis*, pages 1094–1096. Springer Berlin Heidelberg, 2011. ISBN 978-3-642-04898-2. doi: 10.1007/978-3-642-04898-2_455.
- S. M. Kay. *Fundamentals of Statistical Signal Processing: Estimation Theory*. Prentice-Hall, Inc., 1993.
- D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- J. Korner and K. Marton. How to encode the modulo-two sum of binary sources. *IEEE Transactions on Information Theory*, 25(2):219–221, 1979. doi: 10.1109/TIT.1979.1056022.

- A. Krizhevsky. *Learning Multiple Layers of Features from Tiny Images*. PhD thesis, University of Toronto, 2009.
- W. F. Lo, N. Mital, H. Wu, and D. Gündüz. Collaborative semantic communication for edge inference. *IEEE Wireless Communications Letters*, 2023. doi: 10.1109/LWC.2023.3256006.
- N. Mital, E. Özyilkan, A. Garjani, and D. Gündüz. Neural distributed image compression with cross-attention feature alignment. In *IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023. doi: 10.1109/WACV56688.2023.00253.
- M. Nakanoya, S. S. Narasimhan, S. Bhat, A. Anemogianis, A. Datta, S. Katti, S. Chinchali, and M. Pavone. Co-design of communication and machine inference for cloud robotics. *Autonomous Robots*, 2023.
- V. Prabhakaran, D. Tse, and K. Ramachandran. Rate region of the quadratic gaussian CEO problem. In *International Symposium on Information Theory*, 2004. doi: 10.1109/ISIT.2004.1365154.
- J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016. doi: 10.1109/CVPR.2016.91.
- F. Rellich and J. Berkowitz. *Perturbation theory of eigenvalue problems*. CRC Press, 1969.
- M. Salzmann, C. H. Ek, R. Urtasun, and T. Darrell. Factorized orthogonal latent spaces. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, 2010.
- A. Singh, E. Garza, A. Chopra, P. Vepakomma, V. Sharma, and R. Raskar. Decouple-and-sample: Protecting sensitive information in task agnostic data release. In *European Conference on Computer Vision*, 2022. doi: 10.1007/978-3-031-19778-9_29.
- D. Slepian and J. Wolf. Noiseless coding of correlated information sources. *IEEE Transactions on Information Theory*, 1973. doi: 10.1109/TIT.1973.1055037.
- F. Torabi, G. Warnell, and P. Stone. Behavioral cloning from observation. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, 2018. ISBN 9780999241127.
- J. Whang, A. Acharya, H. Kim, and A. G. Dimakis. Neural distributed source coding. *arXiv preprint arXiv:2106.02797*, 2021.
- M. Ye, C. Gong, L. Nie, D. Zhou, A. Klivans, and Q. Liu. Good subnetworks provably exist: Pruning via greedy forward selection. In *International Conference on Machine Learning*, 2020.
- A. Zhan, R. Zhao, L. Pinto, P. Abbeel, and M. Laskin. A framework for efficient robotic manipulation. In *Deep RL Workshop at Advances in Neural Information Processing Systems*, 2022.
- X. Zhang, J. Shao, and J. Zhang. LDMIC: Learning-based distributed multi-view image coding. In *The Eleventh International Conference on Learning Representations*, 2023.
- R. Zhao, U. Topcu, S. Chinchali, and M. Phielipp. Learning sparse control tasks from pixels by latent nearest-neighbor-guided explorations. *arXiv preprint arXiv:2302.14242*, 2023.

Appendix

A. Related Work

Information theoretic perspective: Slepian and Wolf *et al.* are the first to obtain the minimum bandwidth of distributed sources to perfectly reconstruct data (Slepian and Wolf, 1973). However, they use exponentially complex compressors while assuming that the joint distribution of sources is known, which is impractical. In the presence of a task, finding the rate region of two binary sources has remained an open problem, even for modulo-two sum tasks (Korner and Marton, 1979). In terms of imperfectly reconstructing data with neural autoencoders, previous works consider compression of the original data to a fixed dimension (Whang *et al.*, 2021; Diao *et al.*, 2020), while our work focuses on compressing data to any bandwidth with a task model.

Task-aware compression: Real-world data, such as images or audio, are ubiquitous and high-dimensional, while downstream tasks that input the data only utilize certain features for the output. Task-aware compression aims to compress data while maximizing the performance of a downstream task. Previous works analyze linear task (Cheng *et al.*, 2022a), image compression (Ji *et al.*, 2012; Choi and Han, 2020; Nakanoya *et al.*, 2023; Dubois *et al.*, 2021), future prediction (Cheng *et al.*, 2021), and data privacy (han Li *et al.*, 2023; Cheng *et al.*, 2022b), while ours compresses distributed sources under limited bandwidth.

Neural autoencoder: Previous works show the ability of neural autoencoders to generate meaningful and uncorrelated representations. Instead of adding additional loss terms during training like (Bardes *et al.*, 2022; Singh *et al.*, 2022; Chen *et al.*, 2018; Bousmalis *et al.*, 2016; Lo *et al.*, 2023), we use a random projection module to help a neural autoencoder learn uncorrelated and linear-compressible representations. Other works focus on designing new neural architectures for multi-view image compression (Zhang *et al.*, 2023; Mital *et al.*, 2023), while ours focuses on the framework to compress data to different compression levels. We choose autoencoders instead of variational autoencoders (Kingma and Welling, 2013; Higgins *et al.*, 2017) because we focus on the compression of fixed representations rather than generative tasks from latent distributions. Also, autoencoders are more compatible with DPCA than variational autoencoders.

B. Solving DPCA

We now solve the optimization problem in (2). For any given E_1, E_2 (thus, a given Z), we can optimally obtain $M^* = YZ^\top(ZZ^\top)^{-1} = YZ^\top$ by linear regression. Now, we are left to find the optimal encoders E_1, E_2 . First, a preprocessing step removes the correlation part of X_1 from X_2 by subtracting the least-square estimator $\hat{X}_2(X_1)$:

$$\tilde{X}_2 = X_2 - \hat{X}_2(X_1) = X_2 - X_2 X_1^\top (X_1 X_1^\top)^{-1} X_1. \quad (5)$$

The orthogonality principle of least-square estimators (Kay, 1993) ensures that $X_1 \tilde{X}_2^\top = \mathbf{0}_{n_1 \times n_2}$. We decouple the objective in (2a) with respect to E_1, E_2 by the orthogonality principle and (2c):

$$\begin{aligned} \min_{E_1, E_2} \|Y - M^* Z\|_2^2 &= \|Y\|_2^2 - \max_{E_1, E_2} \|M^*\|_2^2 \\ &= \|Y\|_2^2 - \max_{E_1} \|Y_1 X_1^\top E_1^\top\|_2^2 - \max_{E_2} \|Y_2 \tilde{X}_2^\top E_2^\top\|_2^2, \end{aligned} \quad (6)$$

where $Y = \Phi X = [\Phi_1 \Phi_2] [X_1^\top X_2^\top]^\top = Y_1 + Y_2$. We then have two subproblems from (2):

$$\begin{aligned} E_1^* &= \operatorname{argmax}_{E_1} \|\Phi_1 X_1 X_1^\top E_1^\top\|_2^2 \\ \text{s.t.} \quad & E_1 X_1 X_1^\top E_1^\top = \mathbb{I}_{m_1}, \end{aligned} \quad (7)$$

$$\begin{aligned} E_2^* &= \operatorname{argmax}_{E_2} \|\Phi_2 \tilde{X}_2 \tilde{X}_2^\top E_2^\top\|_2^2 \\ \text{s.t.} \quad & E_2 \tilde{X}_2 \tilde{X}_2^\top E_2^\top = \mathbb{I}_{m_2}. \end{aligned} \quad (8)$$

The two subproblems are the canonical correlation analysis (Hardoon *et al.*, 2004), which can be solved by whitening $E_1 X_1, E_2 \tilde{X}_2$ and singular value decomposition (see (Hardoon *et al.*, 2004) for details).

C. Proof of Lemma

C.1. Bounds of DPCA

Lemma (Bounds of DPCA Reconstruction). Given a zero-mean data matrix and its covariance,

$$X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \in \mathbb{R}^{(n_1+n_2) \times N}, XX^\top = \underbrace{\begin{bmatrix} \text{Cov}_{11} & \mathbf{0} \\ \mathbf{0} & \text{Cov}_{22} \end{bmatrix}}_{X_{\text{diag}}} + \underbrace{\begin{bmatrix} \mathbf{0} & \text{Cov}_{12} \\ \text{Cov}_{21} & \mathbf{0} \end{bmatrix}}_{\Delta X},$$

assume that ΔX is relatively smaller than XX^\top , and XX^\top is positive definite with distinct eigenvalues. For PCA's encoding and decoding matrices $E_{\text{PCA}}, D_{\text{PCA}}$ and DPCA's encoding and decoding matrices $E_{\text{DPCA}}, D_{\text{DPCA}}$, the difference of the reconstruction losses is bounded by

$$0 \leq \|X - D_{\text{DPCA}} E_{\text{DPCA}}(X)\|_2^2 - \|X - D_{\text{PCA}} E_{\text{PCA}}(X)\|_2^2 = - \sum_{i=m+1}^{n_1+n_2} \lambda_i e_i^\top \Delta X e_i.$$

where λ_i and e_i are the i -th largest eigenvalue and eigenvector of XX^\top , Tr is the trace function, and m is the dimension of the compression bottleneck.

Proof. The lower bound is intuitive. We know that DPCA cannot outperform PCA since distributed coding cannot outperform joint coding and PCA is the optimal linear encoding. The reconstruction loss of PCA is always not greater than the loss of DPCA, thus the lower bound is 0. Now consider the upper bound:

$$\begin{aligned} & \|X - D_{\text{DPCA}} E_{\text{DPCA}} X\|_2^2 - \|X - D_{\text{PCA}} E_{\text{PCA}} X\|_2^2 \\ &= \text{Tr}(XX^\top + D_{\text{DPCA}} E_{\text{DPCA}} X (D_{\text{DPCA}} E_{\text{DPCA}} X)^\top - 2D_{\text{DPCA}} E_{\text{DPCA}} XX^\top) \\ &\quad - \sum_{i=m+1}^{n_1+n_2} \lambda_i (XX^\top) \\ &= \text{Tr}(X_{\text{diag}} + \Delta X + D_{\text{DPCA}} E_{\text{DPCA}} X (D_{\text{DPCA}} E_{\text{DPCA}} X)^\top - 2D_{\text{DPCA}} E_{\text{DPCA}} XX^\top) \\ &\quad - \sum_{i=m+1}^{n_1+n_2} \lambda_i (XX^\top) \\ &= \text{Tr}(\Delta X + E_{\text{DPCA}}^\top D_{\text{DPCA}}^\top D_{\text{DPCA}} E_{\text{DPCA}} \Delta X - 2D_{\text{DPCA}} E_{\text{DPCA}} \Delta X) \\ &\quad + \sum_{i=m+1}^{n_1+n_2} \lambda_i (X_{\text{diag}}) - \lambda_i (XX^\top) \\ &= \sum_{i=m+1}^{n_1+n_2} \lambda_i (X_{\text{diag}}) - \lambda_i (XX^\top). \end{aligned}$$

Finally, we use the matrix perturbation theory ([Rellich and Berkowitz, 1969](#)) to calculate the first-order approximation of the effect of ΔX on the singular values of X_{diag} . The perturbation theory assumes that the perturbation ΔX is relatively small compared to X_{diag} . Then, we know:

$$\begin{aligned} \|X - D_{\text{DPCA}} E_{\text{DPCA}} X\|_2^2 - \|X - D_{\text{PCA}} E_{\text{PCA}} X\|_2^2 &= \sum_{i=m+1}^{n_1+n_2} \lambda_i (X_{\text{diag}}) - \lambda_i (XX^\top) \\ &\leq \sum_{i=m+1}^{n_1+n_2} \lambda_i - \lambda_i - \lambda_i e_i^\top \Delta X e_i \\ &= - \sum_{i=m+1}^{n_1+n_2} \lambda_i e_i^\top \Delta X e_i. \end{aligned}$$

□

Note that the encoding and decoding matrices of DPCA look like:

$$D_{\text{DPCA}} = \begin{bmatrix} D_1 & \mathbf{0} \\ \mathbf{0} & D_2 \end{bmatrix}, E_{\text{DPCA}} = \begin{bmatrix} E_1 & \mathbf{0} \\ \mathbf{0} & E_2 \end{bmatrix}$$

where E_1, E_2, D_1, D_2 are matrices obtained from each source with DPCA.

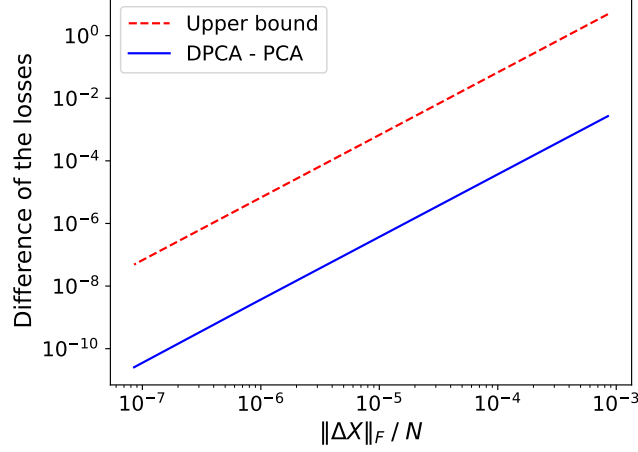


Figure 4: **Bound from Lemma C.1:** The obtained upper bound is always larger than the difference of losses of DPCA and PCA.

We examine the correctness of our bound with random data matrices in Fig. 4. We can see that the gap between DPCA and PCA decreases as the Frobenius norm of ΔX decreases. The upper bound also has the same trend, while it is always larger than the exact value. Note that in Fig. 4, all axes are in log scale.

D. DPCA Module Pseudocode

Algorithm 1 Projection into a random low dimension using DPCA

```

1: Input: A size  $b$  batch of latent representations  $Z_i \in \mathbb{R}^{b \times m_i}$  from each source  $i$ , min and max bandwidth  $m_{\min}, m_{\max}$ 
2: Output: Compressed representation  $Z_i^m$  of each source, reconstructed representation  $\hat{Z}$  for all sources
3: function ENCODE( $Z_i, m_{\min}, m_{\max}$ )
4:   for each source  $i$  do
5:      $\bar{Z}_i \leftarrow Z_i - \text{Mean}(Z_i)$                                      ▷ Normalize representations
6:      $s_i, V_i, H_i \leftarrow \text{SVD}(\bar{Z}_i)$                                ▷ Singular value decomposition
7:   end for
8:    $s, V \leftarrow \text{Cat}(s_i), \text{Cat}(V_i)$                                ▷ Concatenate singular values and vectors
9:    $m \leftarrow \text{Rand}(m_{\min}, m_{\max})$                                ▷ Randomly choose projection dimension
10:   $s^m, V^m \leftarrow \arg \max([s, V], m)$                            ▷ Select the top  $m$  values of  $s$ 
11:  for each source  $i$  do
12:     $V_i^m \leftarrow \{V | V \in V^m, V \in V_i\}$                        ▷ Select  $m$  vectors from sources
13:     $Z_i^m = \bar{Z}_i \times V_i^m$                                            ▷ Project  $Z_i$  to lower dimensions
14:  end for
15:  return  $Z_{\text{low}} \leftarrow \text{Cat}(Z_i^m)$                                ▷ Return Compressed representation
16: end function
17: function DECODE( $Z_i^m$ )
18:  for each source  $i$  do
19:     $\hat{Z}_i \leftarrow Z_i^m \times \text{Cat}(V_i^m)^\top$                            ▷ Decompressed representation
20:     $\hat{Z}_i \leftarrow \hat{Z}_i + \text{Mean}(Z_i)$                                ▷ Denormalize representations
21:  end for
22:  return  $\hat{Z} \leftarrow \text{Cat}(\hat{Z}_i)$                                ▷ Return reconstructed representations
23: end function

```

E. Discussion of Other Experiments

We now describe the other two datasets and the corresponding tasks of our additional experiments:

CIFAR-10 denoising: We first consider a simple task of denoising CIFAR-10 images using two noisy observations of the same image, shown in Fig. 3 (a). Here, the importance of each observation, or view, for the task is simply the noise level. For view 1, we consider an image corrupted with additive white Gaussian noise (AWGN) with a variance of 0.1^2 . And view 2 is highly corrupted by AWGN with a variance of 1. All the images were normalized to $[0, 1]$ before adding the noise. We compressed the noisy observations and passed the reconstructed images through a pre-trained denoising network. We then computed the final peak signal-to-noise ratio (PSNR) with respect to the clean image. Since the noise levels of both views are unequal, the importance of the task is unequal as well. The optimal bandwidth allocation should not be equal, thus showing the advantage of NDPCA. Although view 1 contains more information, not all bandwidth should be allocated to view 1. This problem is called the CEO problem (Berger et al., 1996; Prabhakaran et al., 2004). In fact, even if one view is highly corrupted, we should still leverage that view and never allocate 0 bandwidth to it.

Locate and lift: For the manipulation task, we consider a scenario in which a simulated 6 degrees-of-freedom robotic arm controlled by a reinforcement learning agent inputs two camera views to locate and lift a yellow brick. We call the view from the robotic arm "arm-view" and the one recording the whole desk "side-view", as shown in Fig. 3 (b). The two views are complementary to completing the task, details discussed in Appendix H.3. We trained the agent in a supervised-learning manner. We collected a dataset of observation and action pairs (Zhao et al., 2023) and trained an agent from the dataset. Then, we defined task loss as the L_2 norm of actions from images with and without compression and trained NDPCA to minimize the task loss through the agent. Literature calls this training method "behavior cloning" (Torabi et al., 2018) as it learns from demonstrations. Behavior cloning causes a drop in performance, but this paper only focuses on the performance degradation caused by compression, so we treat the behavior cloning agent with uncompressed views as the upper bound of our method.

The results of the other two experiments are:

Fig. 2(a) shows the results of denoising CIFAR-10 with NDPCA trained at $(m_{\min}, m_{\max}) = (8, 64)$. Although view 1 is more important than view 2, DAE can only equally allocate bandwidth to both sources. NDPCA compresses the data and flexibly allocates bandwidths, as shown in 2(d), where we can see that Z_1 has more bandwidth than Z_2 . NDPCA results in 1.2 dB gain in PSNR compared to DAE when $m = 64$.

Fig. 2(b) shows the results of the locate and lift task with NDPCA trained at $(m_{\min}, m_{\max}) = (8, 48)$. We set the length of an episode as 50 time steps and measure the success rate in 100 episodes. We show the upper bound, a behavior cloning agent without compression, in gray dotted lines. The arm view is more important as it captures the precise location of the brick, and as expected, NDPCA allocates more bandwidth to the arm-view (Z_2), as seen in Fig. 2(e). We see that NDPCA has a 9% higher success rate compared to DAE when $m = 24$.

Comparison of NDPCA with JAE: JAE uses the information from both views simultaneously to capture the best joint embedding for the task. In an ideal scenario, JAE will be the upper bound for the performance and hence easily performs better than DAE across all the experiments. Interestingly, in Fig. 2(b) and (c), we see that NDPCA outperforms not only DAE but also JAE as well. We attribute it to the better representations present in higher-dimension latent space. It turns out that learning a high-dimensional representation and then projecting to a lower dimension space, like NDPCA, is more efficient compared to directly learning a low-dimensional representation, like JAE. This projection from higher dimensional to lower dimensional is similar to pruning large neural networks to identify effective sparse sub-networks. (Frankle and Carbin, 2019; Ye et al., 2020). We also note that Low-Rank Adaptation (LoRA) (Hu et al., 2022) technique for large language models can be thought of as a similar approach.

Limitations: In general, autoencoders are poor at generalizing to out-of-distribution data and the drawback translates to NDPCA as well. When the testing set is noticeably different from the training set, the performance of NDPCA can get noticeably lower. Additionally, during training, DPCA performs the singular value decomposition in the training set. The decomposition operation can become ill-conditioned and unstable if the batch size is too small. An alternative approach could be a parametric low-rank decomposition such as LoRA (Hu et al., 2022) or using adapter networks (Houlsby et al., 2019), although the complexity increases and the compatibility with DPCA remains to be explored.

F. Details of the Datasets

F.1. CIFAR-10 denoising:

We started with the standard CIFAR-10 dataset and normalized the images to $[0, 1]$. Two different views are created by adding different levels of Gaussian noise, $\mathcal{N}(0, 0.1^2)$ and $\mathcal{N}(0, 1)$. The pre-trained task model is created by training a denoising autoencoder that takes both views, concatenates them along the channel dimension, and produces a clean image. The autoencoders need to learn features that are important for this task model.

F.2. Locate and lift:

We collected 20,000 pairs of actions and the corresponding images of both views for our training set. The actions are 4 dimensional, controlling the x, y, z coordinate movements and the gripper of the robotic arm. We randomly cropped the images from 128×128 to 112×112 pixels to make our autoencoder more robust. The expert agent is pre-trained by the same data augmentation as well.

F.3. Airbus detection:

We first cropped all original images of 2560×2560 pixels (Fig. 5) into 224×224 pixels with 28 pixels overlapping between each cropped image. We then eliminated the bounding boxes that are less than 30% left after cropping.



Figure 5: **Original image of airbus detection.** The original images are 2560×2560 pixels, and we cropped them into smaller pieces in 224×224 .

G. Implementation Details

G.1. CIFAR-10 denoising:

For the CIFAR-10 dataset, we used the standard CIFAR-10 dataset and applied different levels of AWGN noise to create two correlated datasets. We used the CIFAR-10 experiments as a proof of concept to try different architectures and loss functions and other techniques to finalize our framework. We choose $\lambda_{\text{task}} = 1$ for the task-aware setting and $\lambda_{\text{rec}} = 1$ for the task-agnostic setting. We run 4 random seeds on NDPCA and all baselines to evaluate the performance.

G.2. Locate and lift:

For the locate and lift experiment, we trained our autoencoder with the same random cropping setting as in Sec. F, which cropped the images from 128×128 to 112×112 pixels. During testing, we randomly initialized the location of the brick and center-cropped the images from 128×128 to 112×112 pixels. We scaled all images to 0 to 1 and ran 5 random seeds on NDPCA and all baselines to evaluate the performance. For the task-aware setting, $\lambda_{\text{task}} = 500$, and $\lambda_{\text{rec}} = 1000$ for the task-agnostic. setting

G.3. Airbus detection:

For the Airbus detection task, we used the original Yolo paper for our object detection model together with the detection loss (Redmon et al., 2016). Our experiments with the latest state-of-the-art Yolo v8 model (Yol, 2023) showed that there is no big difference in the Airbus detection dataset in terms of run time and accuracy. Since the size of the original dataset is not enough to train an object detection model, we used the data augmentation proposed in Yolo v8, mosaic, to increase the size of the dataset. Mosaic randomly crops 4 images and merges them to generate a new image. We used random resized crop, blur, median blur, and CLAHE enhancement during training, each with probability 0.05 by functions in the Albumentations package (Buslaev et al., 2020). We increased the size of the Airbus dataset from 5904 to 21808 with mosiac and trained the Yolo detection model. Finally, we trained our autoencoder with the same dataset, but downsample the images to 112×112 pixels so that the autoencoder is faster to train. For the task-aware setting, $\lambda_{\text{task}} = 0.1$, and $\lambda_{\text{rec}} = 0.5$ for the task-agnostic setting. We run 2 random seeds on NDPCA and all baselines to evaluate the performance.

G.4. Neural Autoencoder Architecture and Hyperparameters

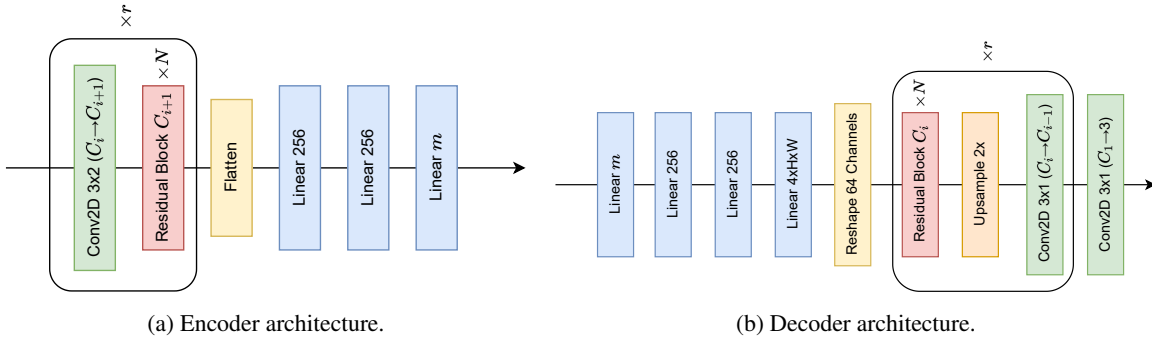


Figure 6: **ResNet Autoencoder**: The encoder processes inputs through r convolution layers and $r \times N$ residual blocks, followed by 3 fully connected layers with ReLU activation. The decoder processes latent representations in the reverse order from the encoder with $2 \times$ upsamplings.

We used the ResNet encoder shown in Fig. 6a and the decoder in Fig. 6b for all experiments. We used different numbers of filters and numbers of residual blocks for our experiments, shown as C and r . We denote m as the number of latent dimensions. The numbers of filters are $C_1 = 32$, $C_2 = 64$, $C_3 = 128$, $C_1 = 8$, $C_2 = 16$, $C_3 = 32$, $C_4 = 64$, and $C_1 = 16$, $C_2 = 32$, $C_3 = 64$, $C_4 = 128$, and the numbers of residual blocks are $r = 0$, $r = 1$, $r = 1$ for CIFAR-10 denoising, locate and lift, and Airbus detection. For CIFAR-10 denoising, we use the Adam optimizer with a learning rate of 0.0002, and for the other two experiments, we use the Adam optimizer with a learning rate of 0.0001. For the sake of training speed, when training DAE and JAE, we first trained a large network with m_{max} with each random seed. Then, we fixed the network parameters and trained concatenate 3 fully connected layers on each encoder and decoder network to compress and decompress the data to smaller m .

G.5. Balancing Task-aware and Task-agnostic Loss

NPDCa has a loss function consisting of 2 terms, as shown in (4):

$$\mathcal{L}_{\text{tot}} = \lambda_{\text{task}} \underbrace{\|\hat{Y} - Y\|_F^2}_{\text{task loss}} + \lambda_{\text{rec}} \underbrace{\left(\|\hat{X}_1 - X_1\|_F^2 + \|\hat{X}_2 - X_2\|_F^2 + \dots + \|\hat{X}_K - X_K\|_F^2 \right)}_{\text{reconstruction loss}}. \quad (4 \text{ revisited})$$

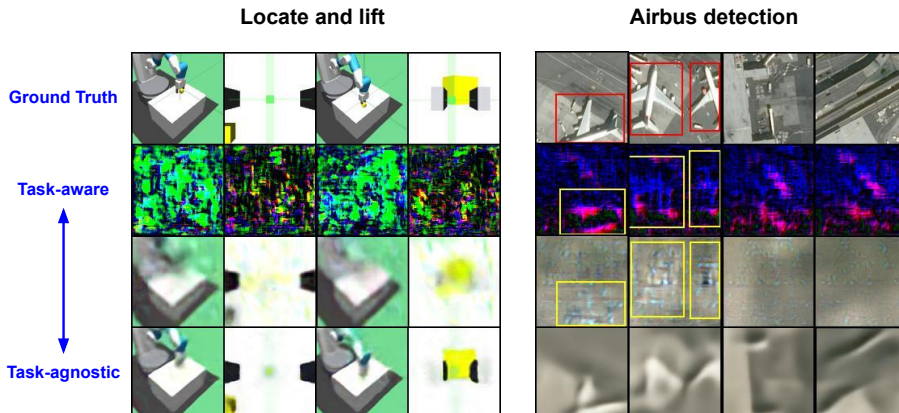


Figure 7: **Weighted task-loss:** Weighted task-aware images faintly reconstruct the original images while restoring task-relevant features with high-frequency noise. In Airbus detection, location of Airbuses is captured with shiny high-frequency pixels in row 3.

Previous work (Nakanoya et al., 2023) tested cases of (4), such as task-aware when $\lambda_{\text{task}} > 0, \lambda_{\text{rec}} = 0$, and task-agnostic when $\lambda_{\text{task}} = 0, \lambda_{\text{rec}} > 0$. Of course, one can use different weighted sums of the 2 terms in (4), which we call weighted task-aware. We show the resulting reconstructed image in Fig. 7, whose weights are a mixture of half of the two other methods. Weighted task-aware images have both blurry reconstructions of the original images and task-relevant features. Unsurprisingly, the task loss and the reconstructed loss of weighted task-aware images are between pure task-aware and task-agnostic, that is, we can use the weights in the loss function to trade off compressing human perception features against task-relevant features. Interestingly, we can see that the task-aware images look similar to the images without Airbuses (last 2 columns), and when there are Airbuses, the task-aware images look different. It means that the features of no Airbuses are pretty much the same in the latent space, thus resulting in similar images in pixel space. Hence we can conclude that task-aware features are not random noise, they are meaningful features only to the task model but not to our eyes.

G.6. Storage and Training Complexity

Model	CIFAR-10		Locate and lift		Airbus detection	
	Storage (MB)	Train (hr)	Storage (MB)	Train (hr)	Storage (MB)	Train (hr)
NDPCA	8.3	0.25	16.4	5.0	33.0	13.0
DAE	5×8.4	5×0.21	4×16.3	4×5.0	4×22.5	4×11.5
JAE	5×10.2	5×0.22	4×11.4	4×3.5	4×32.9	4×10.5

Table 1: **Storage and training complexity:** NDPCA has slightly more storage and training overload than other models for a single bandwidth but can operate across different bandwidths. We multiply the number of bandwidths tested in Fig. 2 to the storage size and training time of DAE and JAE as they require different models for different compression levels.

One key feature of NDPCA is that it only needs one model to operate in different bandwidths. Therefore, we only need to train and store one model at the edge devices and the central node. We compare the complexity of storage and training in Table 1. Although NDPCA has a larger storage size and longer training time than other models, it can operate across different bandwidths. According to Table 1, if all models operate in more than 1 bandwidths, NDPCA saves more storage and training overload because other models have more than 50% of NDPCA’s overload. For CIFAR-10 denoising, we tested the training time on an RTX 4090, and for the locate and lift and Airbus detection experiments, we tested the training time on an NVIDIA RTX A5000.

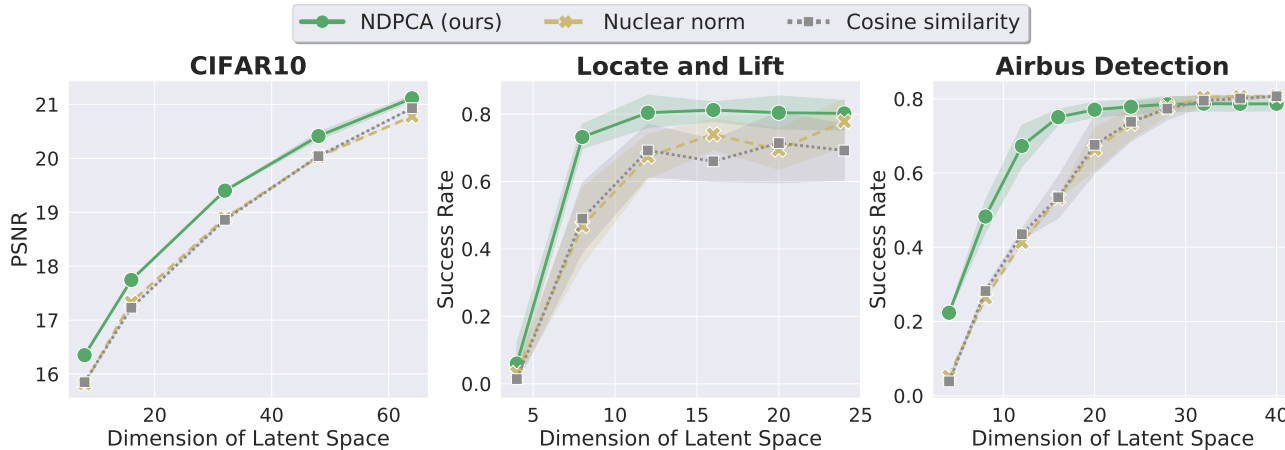


Figure 8: **Ablation study of the nuclear norm and cosine similarity:** Adding the nuclear norm or cosine similarity to the loss function does not improve the performance of the model when compressing latent representations to lower dimensions.

H. Ablation Study

H.1. Cosine similarity and nuclear norm

In Fig. 8, we show that adding nuclear norm or cosine similarity in the training loss (4) does not help the model perform when we use DPCA to project latent representations into lower dimensions. We compared our proposed NDPCA with the DPCA module against NDPCA without the DPCA module but with the penalization of the nuclear norm and cosine similarity added. The weights of all the additional terms are 0.1. From Fig. 8, we conclude that the DPCA module can increase the performance better than the other two.

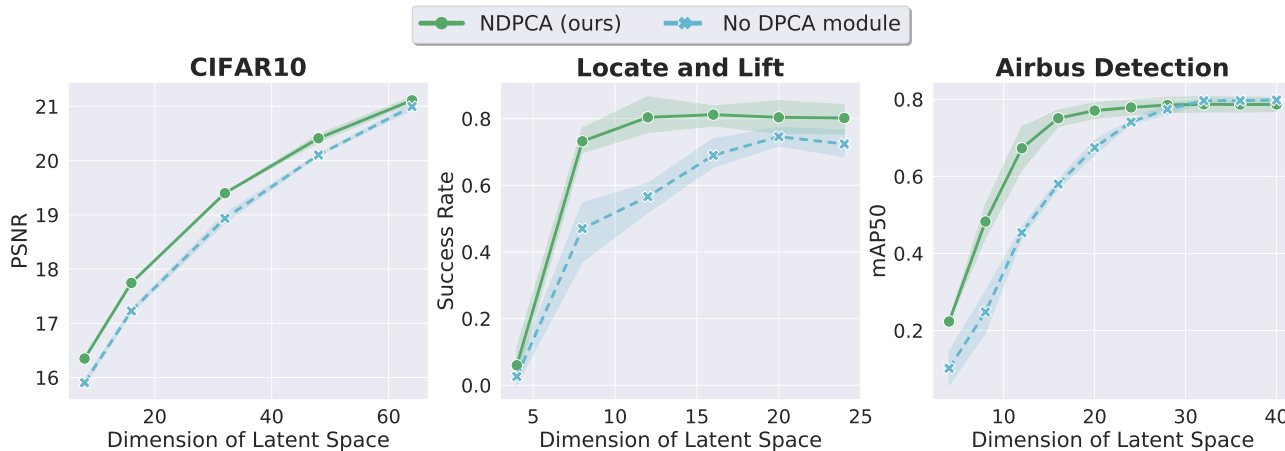


Figure 9: **Ablation study of DPCA module:** The proposed DPCA module effectively increases the performance in lower bandwidths, while achieving the same performance at larger bandwidths.

H.2. DPCA module

In Fig. 9, we show that the proposed DPCA module can help the neural autoencoder learn linear compressible representations, as described in Sec. 4. We see that with the DPCA module, NDPCA can increase the performance in lower bandwidths, while saturating at the performance close to the model without the module. We conclude that with the DPCA module, NDPCA learns to generate low-rank representations, so the performance is better in lower bandwidths. However, when the bandwidth is higher, the bandwidth can almost fully restore the representations, so the two methods perform similarly.

H.3. Single view performance of locate and lift

In the locate and lift experiments, the reinforcement learning agent leverages information from both views as input to manipulate. Here, we detail why the 2 views are complementary to accomplish the task. The success rate of an agent is 76% with only the arm-view and 45% with the side-view. When combining both, the success rate is 83%. The reason why the views are complementary is that the side-view provides global information on the position of the arm and the brick, but sometimes the brick is hidden behind the arm. The arm-view captures detailed information from a narrow view of the desk. Once the arm-view captures the brick, it is straightforward to move toward it and lift it. The arm view is more important because with only the arm-view, the agent can randomly explore the brick, but with only the side-view, the brick might be vague to see and thus harder to lift. Of course, with both views, the robotic arm can easily move toward the vague position of the brick and use arm-view to lift it.