

# Co-Design of Communication and Machine Inference for Cloud Robotics

Manabu Nakanoya · Sai Shankar Narasimhan\* · Sharachchandra Bhat\* · Alexandros Anemogiannis · Akul Datta · Sachin Katti · Sandeep Chinchali · Marco Pavone

Received: date / Accepted: date

**Abstract** Today, even the most compute-and-power constrained robots can measure complex, high data-rate video and LIDAR sensory streams. Often, such robots, ranging from low-power drones to space and subterranean rovers, need to transmit high-bitrate sensory data to a remote compute server if they are uncertain or cannot scalably run complex perception or mapping tasks locally. However, today’s representations for sensory data are mostly designed for *human, not robotic*, perception and thus often waste precious compute or wireless network resources to transmit unimportant parts of a scene that are unnecessary for a high-level robotic task. This paper presents an algorithm to learn *task-relevant* representations of sensory data that are co-designed with a pre-trained robotic perception model’s ultimate objective. Our algorithm aggressively compresses robotic sensory data by up to  $11 \times$  more than competing methods. Further, it achieves high accuracy and robust generalization on diverse tasks including Mars terrain classification with low-power deep learning accelerators, neural motion planning, and environmental timeseries classification.

---

Manabu Nakanoya  
 NEC Corporation, Tokyo, Japan  
 E-mail: {nakanoya}@nec.com

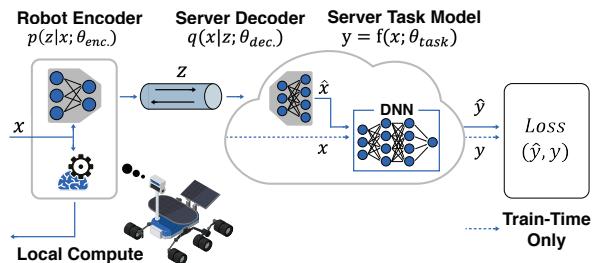
Sai Shankar Narasimhan, Sharachchandra Bhat, and Sandeep Chinchali  
 The University of Texas at Austin,  
 Department of Electrical and Computer Engineering  
 E-mail: {nsaishankar, sharachchandra, sandeepc}@utexas.edu

Sachin Katti and Marco Pavone  
 Stanford University, Department of Computer Science  
 E-mail: {skatti, pavone}@stanford.edu

\* refers to equal contribution.

## 1 Introduction

Imagine a future Mars or subterranean rover that captures high-bitrate video and LIDAR sensory streams as it charts uncertain terrain, some of which it cannot classify locally. How should these robots represent, compress, and transmit their rich sensory data over bandwidth-limited wireless networks, especially if the intended audience is often a compute-intensive, remote machine learning model, not necessarily a human viewer? Indeed, even a single RGB-D (depth) camera stream produces upwards of 45 Megabytes/second of data [30], while the deep-space network only has a communication bandwidth of 0.5-4 Megabits/second [2].<sup>1</sup>



**Fig. 1: Task-Relevant Communication for Perception:** A compute-limited robot *learns* how to compress sensory input  $x$ , transmit salient features  $z$ , and decode the input  $\hat{x}$  so that it can directly leverage a pre-trained, potentially “off-the-shelf” task module  $f(\cdot; \theta_{\text{task}})$  at a central server. By learning an encoder and decoder (gray) within the context of a pre-trained task module’s goal, we only transmit minimal, salient information. In the figure, only components in gray (the encoder and decoder) are learned, while the task network  $f$  has pre-trained, fixed parameters. Original sensory input  $x$  and task output  $y$  (dashed lines) are only used for training.

More broadly, today’s representations for sensory data mostly optimize for human, not robotic, perception and thus

<sup>1</sup> NASA estimate for Earth to Mars Reconnaissance Orbiter link [2].

try to faithfully represent every pixel or point-cloud in a scene [5, 10]. Ideally, however, resource-constrained robots should represent only salient parts of sensory streams for remote perception and planning tasks to reduce the computational cost of encoding, storing, and transmitting sensory data. This paper presents a general algorithmic framework to *learn* such concise, task-relevant representations. We note that several recent works have used specialized deep neural networks (DNNs) to improve LIDAR or JPEG image compression [40, 26, 43], motivated by the observation that standard compression schemes do not emphasize features that are most sensitive for accurate *machine* perception [10, 5, 3]. However, the key novelty of our work is to use a general-purpose, pre-trained task module to *guide* representation learning, which allows us to easily generalize to multiple sensor modalities unlike prior specialized point-solutions.

Specifically, as shown in Fig. 1, our design provisions a pre-trained, differentiable task module at a central server, which allows multiple robots to share an upfront cost of training the model and benefit from centralized model updates. Our key technical insight is to *co-design* a minimal sensory representation that is tightly coupled with the task module’s objective. Henceforth, co-design means that the pre-trained task network parameters are fixed and the task objective guides what salient parts of a sensory stream to encode. By also learning a task-relevant decoder, we enable robots to leverage a variety of pre-trained, publicly-available task modules without modifying their required input dimensions. Our approach is complementary to advances in on-robot computation and instead pertains to compute-intensive, potentially collaborative, tasks that require remote assistance. Indeed, our work is partly inspired by recent proposals to place powerful servers in deep space to assist future rovers and space telescopes with compute-intensive machine learning tasks, data processing, and storage before relaying selected data back to earth [13, 42].

**Literature Review:** Recent work has applied information bottleneck theory [38] to build controllers that focus on actionable, task-relevant visual inputs for robust, generalizable navigation and grasping policies [32, 31, 35]. In contrast, we introduce a novel algorithm for co-designing communication and machine perception, which uses pre-trained task modules to learn salient, efficiently-computable representations. Our co-design approach differs from specialized solutions that use DNNs for JPEG image compression [26, 43] and instead uses diverse, off-the-shelf task modules to learn representations for images, obstacle point clouds, and environmental sensor timeseries.

Our work is also related to cloud robotics [20, 23], where robots use remote servers to augment their grasping [25], object recognition [19], and mapping capabilities [28]. Prior work balances cloud computing accuracy with communication delay by learning when to query the cloud [6, 37]. Rather

than address *when* to communicate, we instead address how to represent data for task-centric communication. Our work is also related to variational and multi-task autoencoders [21, 14], which compress inputs to minimize regularized reconstruction loss, often with domain-specific output layers in the multi-task case. In contrast, by instead focusing on task loss, we compress sensory data 11× more than standard autoencoders. [We also note previous contributions in networked robotics like \[18\]](#), which considers multi-robot systems communicating with each other over a bandwidth constrained network. [18] introduces the dynamic information flow metric that provides a trade-off between communication resources required to transmit rich sensor observations and the computation cost of processing them. Our approach differs from this solution as we focus on compressing sensor observations in a task-centric manner.

Finally, our work is complementary to methods that compress large DNNs to run on compute-limited robots, using methods like weight pruning or knowledge distillation [4]. Despite such advances, state-of-the-art DNNs, such as transformers for NLP (e.g., BERT [8]) or vision are still infeasible to run on low-power, compute-limited robots. Thus, our work addresses a viable alternative to transmit salient features for remote, compute-intensive machine learning tasks.

**Statement of Contributions:** In light of prior work, our contributions are three-fold. First, we introduce a novel formulation that co-designs a distributed encoder/decoder with a fixed task network. Second, we develop a novel algorithm that aggressively compresses sensory data for diverse robotic tasks. Third, we show how to flexibly allocate computation between a compute-limited robot and server, which renders our approach compatible with low-power DNN accelerators like the Google Edge Tensor Processing Unit (TPU) [1].

**Organization:** This paper is organized as follows. In Sec. 2, we introduce a novel, general problem of co-designing sensory representations for robotic perception. To address this problem, Sec. 3 presents our co-design algorithm and highlights significant compression gains for an illustrative example of linear systems. Then, Sec. 4 evaluates our algorithm on diverse perception tasks with compute-efficient DNNs. Sec. 5 concludes with future directions.

An earlier version of this paper was presented at the Robotics: Science and Systems Conference in 2021 [29]. In this extended version, we show that the compression benefits of our co-design algorithm generalize well to image datasets like CIFAR-10 (Sec. 4.6). In addition, we perform a detailed analysis on the effects of our co-design algorithm on compressed sensory data. In Sec. 4.8 we show that our co-design algorithm clusters input sensory data in the compressed space much more distinctly than other methods. We focus on the human-interpretability aspect of task-relevant compression in Sec. 4.7. We propose a method that utilizes an adversarially trained robust task network to incentivize

more human-interpretable reconstruction while achieving higher compression and task performance than competing benchmarks.

## 2 Problem Statement

We now introduce the core compute modules for information flow between a robot and a central server (Fig. 1), and then formalize our problem statement. First, a robot measures a sensory input  $x \in \mathbb{R}^n$ , such as an image or LIDAR point cloud. Without loss of generality,  $x$  could represent a single sensory sample a robot wishes to transmit or a window  $w$  of correlated samples measured from time  $t-w$  to  $t$ , such as a segment of video, denoted by  $\mathbf{x} = x^{t-w:t}$ .

**Robot Encoder:** The encoder maps raw input  $x$  to a concise representation  $z \in \mathbb{R}^Z$ , denoted by  $z = p(z|x; \theta_{\text{enc.}})$ , where  $\theta_{\text{enc.}}$  are encoder model parameters, such as learned DNN weights. Henceforth,  $z$  is referred to as a bottleneck representation, since it is compressed via an information bottleneck, such as a DNN hidden layer, to a size  $Z \ll n$ .

**Server Decoder:** The encoder transmits bottleneck representation  $z$  over a wireless link to a differentiable decoder at a central server, which generates a reconstructed estimate of the raw sensory input  $\hat{x} = q(x|z; \theta_{\text{dec.}})$ . The parameters of a decoder model, such as a DNN, are denoted by  $\theta_{\text{dec.}}$ .

**Server Task Network:** Finally, the decoded input  $\hat{x}$  is *directly passed* through a pre-trained, differentiable task module, which serves as the key distinction of our work from a standard autoencoder. An example of a task module could be a DNN object detector that predicts object locations and classes  $\hat{y} = f(\hat{x}; \theta_{\text{task}})$  using model parameters  $\theta_{\text{task}}$ . Crucially, by decoding to the original input dimension  $n$ , our approach enables robots to directly utilize a plethora of pre-trained, publicly-available task modules that expect an input dimension  $n$  without re-training on custom bottleneck representations  $z$ . In practice, the encoder, decoder, and task module can all be DNNs. We assume sensory input  $x$  and its label  $y$  are drawn from a domain-specific distribution  $\mathcal{D}$ , such as the space of all Mars terrain images and labels.

### 2.1 Task and Reconstruction Optimization Objectives

Our principal objective is to minimize the *task loss* defined as  $\mathcal{L}_{\text{task}}(y, \hat{y}; \theta_{\text{task}})$ , which compares the resultant task outputs  $y = f(x; \theta_{\text{task}})$  using original input  $x$  and predicted outputs  $\hat{y} = f(\hat{x}; \theta_{\text{task}})$  using a decoded input  $\hat{x}$ . Our key technical insight is that robotic perception tasks can often tolerate considerable distortion in the input estimate  $\hat{x}$ , such as omitting irrelevant parts of an image/map for classification/planning, as long as the downstream task module  $f(\cdot; \theta_{\text{task}})$  can still achieve its goal. Thus, decoded input  $\hat{x}$  should only repre-

sent task-relevant features, which enables highly-compressed representations  $z$  to be transmitted over a wireless link.

Optionally, a roboticist might want to also minimize reconstruction loss, in cases where decoded inputs  $\hat{x}$  should be moderately human-interpretable to view or debug portions of an image. For such scenarios, the reconstruction loss function  $\mathcal{L}_{\text{recon.}}(x, \hat{x}; \theta_{\text{enc.}}, \theta_{\text{dec.}})$  incentivizes faithful reconstruction of a scene, which, in practice, could be the standard variational autoencoder regularized loss [21]. We provide a general, flexible co-design framework that allows a roboticist to optimize a weighted combination of task and reconstruction loss, with reconstruction loss weight  $\lambda \geq 0$ :

$$\begin{aligned} \mathcal{L}_{\text{weight}}(x, \hat{x}, y, \hat{y}; \theta_{\text{task}}, \theta_{\text{enc.}}, \theta_{\text{dec.}}) &= \\ \mathcal{L}_{\text{task}}(y, \hat{y}; \theta_{\text{task}}) + \lambda \mathcal{L}_{\text{recon.}}(x, \hat{x}; \theta_{\text{enc.}}, \theta_{\text{dec.}}). \end{aligned} \quad (1)$$

Our experiments evaluate the scenario of strictly optimizing for task loss ( $\lambda = 0$ ) as well as various values of  $\lambda > 0$ , which introduces a regularization term to incentivize highly-compressed representations that still yield human-interpretable reconstructions. Having defined our optimization objective, we now formalize our problem statement.

### 2.2 Problem Statement

#### Problem 1 (Sensory Co-design for Machine Perception)

Given a differentiable task module  $f(\cdot; \theta_{\text{task}})$  with fixed, pre-trained parameters  $\theta_{\text{task}}$ , fixed bottleneck dimension  $Z$ , and reconstruction loss weight  $\lambda \geq 0$ , find robot encoder and server decoder parameters  $\theta_{\text{enc.}}$  and  $\theta_{\text{dec.}}$  that minimize weighted loss (Eq. 1) over data distribution  $\mathcal{D}$ :

$$\theta_{\text{enc.}}^*, \theta_{\text{dec.}}^* = \underset{\theta_{\text{enc.}}, \theta_{\text{dec.}}}{\operatorname{argmin}} \mathbb{E}_{(x, y) \sim \mathcal{D}} \mathcal{L}_{\text{weight}}(x, \hat{x}, y, \hat{y}; \theta_{\text{task}}, \theta_{\text{enc.}}, \theta_{\text{dec.}}),$$

where  $\hat{x} = q(p(x; \theta_{\text{enc.}}); \theta_{\text{dec.}})$  and  $\hat{y} = f(\hat{x}; \theta_{\text{task}})$ .

Prob. 1 is widely applicable to resource-limited robots, such as space and mining rovers, that can unlock large benefits of remote computation if they send *task-relevant* information over a bandwidth-limited network. The key novelty of our formulation is that we leverage knowledge of a task objective, through  $f(\cdot; \theta_{\text{task}})$ , to guide task-relevant representation learning of encoder/decoder parameters  $\theta_{\text{enc.}}, \theta_{\text{dec.}}$ . Our formulation is related to multi-task learning for vision, which often considers semantically-similar tasks, such as segmentation and object detection [7]. In contrast, however, our weighted objective has two semantically-different (and often competing) goals of learning a sparse set of features for machine perception while also regularizing for pixel-wise reconstruction for human-interpretability. Further, we enable a roboticist to set bottleneck size  $Z$  in a precise manner based on a network’s maximum allowable data-rate. For example, the data-rate could be the size of  $Z$  32-bit floating point values sent at a certain communication frequency.

### 2.3 Illustrative Example: Linear Systems Setting

To illustrate the compression benefits of optimizing for task loss in Prob. 1, we consider a toy example where the encoder, decoder, and task module are matrices. Specifically, consider a simple robotic sensor network where a central server must estimate a function  $y = Kx$  of a potentially large sensor measurement  $x$ , *without* necessarily sending the full input  $x$  over a wireless link. First, robot encoder matrix  $A = \theta_{\text{enc}} \in \mathbb{R}^{Z \times n}$  generates encoding  $z = p(x; \theta_{\text{enc}}) = Ax$ . After  $z$  is sent over a wireless link, it is decoded to  $\hat{x} = q(z; \theta_{\text{dec}}) = Bz$  by linear decoder  $B = \theta_{\text{dec}} \in \mathbb{R}^{n \times Z}$ . Finally, task matrix  $K \in \mathbb{R}^{m \times n}$  generates the output of the linear estimation problem  $y \in \mathbb{R}^m$ . The task loss penalizes error in the linear estimate  $y = Kx$ , and the reconstruction loss optionally penalizes error in the sensor estimate  $\hat{x}$  in case elements of  $\hat{x}$  need to be sanity-checked. Both the task and reconstruction loss are quadratic, yielding weighted loss:

$$\mathcal{L}_{\text{weight}} = \underbrace{\|(Kx - KB\hat{x})\|_2^2}_{\text{task loss}} + \lambda \underbrace{\|(x - B\hat{x})\|_2^2}_{\text{recon.loss}}. \quad (2)$$

To show the full benefits of task-based compression, we now provide an analytical solution for Prob. 1 for the special case when  $\lambda = 0$  (pure task loss objective).

**Theorem 1 (Linear Task-Aware Compression)** Consider task matrix  $K \in \mathbb{R}^{m \times n}$  with rank  $r$  and compact singular value decomposition (SVD)  $K = U\Sigma V^\top$ , where  $U \in \mathbb{R}^{m \times r}$  and  $V^\top \in \mathbb{R}^{r \times n}$  are semi-unitary. Then, for bottleneck dimension  $r$ , setting  $A = V^\top$  and  $B = V$  solves Prob. 1 with zero task loss  $\|(Kx - KB\hat{x})\|_2^2$  (Eq. 2 with  $\lambda = 0$ ) for any  $x \in \mathbb{R}^n$ . Further, there are no encoder and decoder matrices with bottleneck  $Z < r$  that achieve zero task loss.

*Proof* The task loss is zero when, for any  $x \in \mathbb{R}^n$ ,  $y = Kx = \hat{y} = KB\hat{x}$ . This is achieved when  $B = V$  and  $A = V^\top$ :

$$\begin{aligned} \hat{y} &= KB\hat{x} = (\underbrace{U\Sigma V^\top}_{\text{SVD}(K)}) \times \underbrace{V}_{B} \times \underbrace{V^\top}_{A} \times x \\ &= U\Sigma V^\top x = Kx = y, \end{aligned}$$

where  $V^\top V = I_{r \times r}$  since  $V$  is semi-unitary. We now show there are no other solutions with bottleneck  $Z < r$ . For the sake of contradiction, suppose there exist  $\tilde{A} \in \mathbb{R}^{Z \times n}$ ,  $\tilde{B} \in \mathbb{R}^{n \times Z}$  with  $Z < r$  that achieve zero task loss. Then:

$$\begin{aligned} \text{rank}(\tilde{B}\tilde{A}) &\leq \min(\text{rank}(\tilde{A}), \text{rank}(\tilde{B})) \leq Z, \text{ and} \\ \text{rank}(K\tilde{B}\tilde{A}) &\leq \min(\text{rank}(K), \text{rank}(\tilde{B}\tilde{A})) \leq Z. \end{aligned}$$

However, in order for  $\tilde{A}, \tilde{B}$  to achieve zero task loss, we must have  $Kx = K\tilde{B}\tilde{A}x$  for any  $x \in \mathbb{R}^n$ , meaning we must have  $K = K\tilde{B}\tilde{A}$ . However, this is a contradiction since  $\text{rank}(K\tilde{B}\tilde{A}) \leq Z < r = \text{rank}(K)$ . Thus, we conclude there are no solutions that achieve zero task loss for  $Z < r$ .

**Compression benefits:** Our key insight is that, by co-designing with fixed task matrix  $K$ , we achieve zero task loss with a compression gain of  $\frac{n}{Z}$  compared to sending full input  $x$ , where  $Z$  is the bottleneck and  $Z = \text{rank}(K)$ . However, we do not necessarily have zero reconstruction loss since matrix  $V^\top$  is only semi-unitary, meaning it is not required for  $VV^\top = BA$  to be an identity matrix. In fact, if we want zero reconstruction loss, then we have  $\hat{x} = BAx = x$  only when  $B = A^{-1}$ , requiring  $B$  and  $A$  to be square  $n \times n$  matrices, yielding no compression gain. Crucially, our solution differs from naively just computing  $y = Kx$  directly at the robot, since we only transmit a small representation  $z = V^\top x$  and use the modular task matrix  $K$  at the server. Finally, we note that the choice of  $A = V^\top$  and  $B = V$  is not unique for bottleneck  $r$ . For example, we can scale the encoder to  $aA$  and decoder to  $bB$  where  $ab = 1$ .

### 2.4 Illustrative Example: Weighted Linear Setting

We now generalize our previous result for the *weighted* linear setting with  $\lambda > 0$ . The key highlight of our solution is that the optimal linear encoder/decoder to minimize weighted task loss are simply a solution to a low-rank approximation problem. To establish this result, we first concatenate  $N$  examples drawn from training dataset  $(x, y) \sim \mathcal{D}$  into a sample matrix  $\mathbf{X} \in \mathbb{R}^{n \times N}$ . Given the dataset  $\mathbf{X}$  and a bottleneck constraint  $Z$ , minimizing the weighted task loss (Eq. 2) in the linear setting amounts to solving:

$$\begin{aligned} \underset{A, B}{\operatorname{argmin}} \sum_{i=1}^N (\hat{\mathbf{X}}_i - \mathbf{X}_i)^\top (K^\top K + \lambda I) (\hat{\mathbf{X}}_i - \mathbf{X}_i) \quad &\text{where} \\ \hat{\mathbf{X}} &= BAX, \text{rank}(B) \leq Z \text{ and } \text{rank}(A) \leq Z. \end{aligned} \quad (3)$$

In the above problem, the  $i$ -th column of  $\mathbf{X}$  and  $\hat{\mathbf{X}}$  are given by  $\mathbf{X}_i$  and  $\hat{\mathbf{X}}_i$  respectively. The following theorem solves Problem 1 for the weighted linear setting (Eq. 2) by reducing it to a canonical low-rank approximation problem.

**Theorem 2 (Linear Weighted Compression)** The linear setting of Problem 1 is a canonical low-rank approximation problem with an analytical solution for an optimal encoder matrix  $A$  and decoder matrix  $B$ .

*Proof (Sketch)* We first note that the real matrix  $K^\top K + \lambda I$  is positive semi-definite, which admits an eigen-decomposition denoted by  $Y\Lambda Y^\top$ , where  $Y$  is orthogonal. Denoting the Frobenius norm of a matrix as  $\|\cdot\|_F$ , we can write the weighted task loss for Problem 1 in the linear setting as:

$$\begin{aligned} &\sum_{i=1}^N (\hat{\mathbf{X}}_i - \mathbf{X}_i)^\top (K^\top K + \lambda I) (\hat{\mathbf{X}}_i - \mathbf{X}_i) \\ &= \sum_{i=1}^N (\hat{\mathbf{X}}_i - \mathbf{X}_i)^\top (Y\Lambda Y^\top) (\hat{\mathbf{X}}_i - \mathbf{X}_i) \\ &= \|\Lambda^{\frac{1}{2}} Y^\top \hat{\mathbf{X}} - \Lambda^{\frac{1}{2}} Y^\top \mathbf{X}\|_F^2. \end{aligned}$$

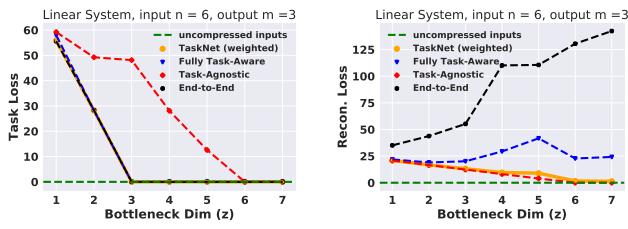


Fig. 2: **Task-aware compression for linear systems.** When the encoder, decoder, and task module  $K$  are matrices, we achieve zero task loss by only sending a representation of size  $\text{rank}(K)$  compared to a measurement of size  $n$ . Our co-design scheme (orange, blue) yields lower task loss for much smaller  $z$  compared to task-agnostic methods (red). Further, we achieve better reconstruction loss and more stable training than an end-to-end method (black), which does not guide learning by fixing task module  $K$ .

Thus, the solution to Problem 1 with a bottleneck  $Z$  in the linear setting can be written in the following low-rank approximation form:

$$\underset{\substack{A, B \\ \text{approximation}}}{\text{argmin}} \| \underbrace{\Lambda^{\frac{1}{2}} Y^\top B A X}_{\text{original}} - \underbrace{\Lambda^{\frac{1}{2}} Y^\top X}_{\text{original}} \|_F^2 \quad \text{where} \\ \text{rank}(B) \leq Z \text{ and } \text{rank}(A) \leq Z. \quad (4)$$

Thus, as stipulated by the Eckhart-Young theorem, the solution to Eq. 4 is the rank  $Z$  truncated singular value decomposition (SVD) of matrix  $\Lambda^{\frac{1}{2}} Y^\top X$ , denoted by  $U \Sigma V^\top$ . For the SVD,  $U \in \mathbb{R}^{n \times Z}$  is semi-orthogonal,  $V \in \mathbb{R}^{m \times Z}$  is semi-orthogonal, and  $\Sigma \in \mathbb{R}^{Z \times Z}$  represents a diagonal matrix of singular values. Thus, an encoder of  $A = U^\top \Lambda^{\frac{1}{2}} Y^\top$  and decoder of  $B = (\Lambda^{\frac{1}{2}} Y^\top)^{-1} U$  are solutions since:

$$\underbrace{\Lambda^{\frac{1}{2}} Y^\top B A X}_{\text{approximation}} = \Lambda^{\frac{1}{2}} Y^\top \underbrace{(\Lambda^{\frac{1}{2}} Y^\top)^{-1} U U^\top \Lambda^{\frac{1}{2}} Y^\top X}_{B} = \underbrace{U(U^\top \Lambda^{\frac{1}{2}} Y^\top X)}_{\text{optimal rank } Z \text{ approximation}} = \underbrace{U \Sigma V^\top}_{\text{optimal rank } Z \text{ approximation}}.$$

**Weighted Linear Results:** For the weighted  $\lambda > 0$  case, we solve Prob. 1 using the SVD as in Theorem 2. Fig. 2 shows results for a toy linear problem with input and output dimensions  $n = 6$  and  $m = 3$ , where we achieve optimal task performance with a  $2\times$  compression gain. Our toy example shows that, even for linear systems, we achieve large compression gains by co-designing with a task objective, which mirror our subsequent results with complex DNNs. While we note our theory only applies to linear systems, such a theoretical underpinning is informative to understand our results for complex DNNs.

### 3 Task-Aware Co-design Algorithm

In general, analytically solving Problem 1 is challenging, especially when the task network and encoder/decoder are

DNNs. We now present Algorithm 1, which learns an encoder/decoder that approximately solves Problem 1. Importantly, our algorithm yields experimental results that match analytical calculations for simple linear settings and also scales to deep learning tasks, presented in Sec. 4.

Alg. 1 takes as input a bottleneck dimension  $Z$ , set based on communication data-rate limits, and user-desired reconstruction weight  $\lambda$ , on line 1. Then, it randomly initializes the encoder/decoder parameters on line 2, which are learned during training *unlike* task parameters  $\theta_{\text{task}}$  (line 3).

The main loop repeats for  $T$  learning rounds, where we obtain training samples of batch-size  $B$  from the labeled dataset  $\mathbf{D} = \{(x_i, y_i)\}_{i=1}^N$ , where  $N$  is the dataset size (line 5), encode/decode sensory data (lines 6-7), and finally, invoke the task module on line 8. The key step for co-design is line 9, where we backpropagate weighted loss  $\mathcal{L}_{\text{weight}}$  (Eq. 1) to update parameters  $\theta_{\text{enc.}}, \theta_{\text{dec.}}$ , while keeping the task parameters  $\theta_{\text{task}}$  *fixed*. Importantly, Alg. 1 can be run *offline* given a training dataset and desired bottleneck  $Z$  to avoid passing large gradients over a wireless link. Then, a robot can deploy the trained encoder and decoder, and periodically improve them with more field data.

---

#### Algorithm 1: Task-Relevant Compression Co-design

---

```

1 Set bottleneck dimension  $Z$ , reconstruction weight  $\lambda$ 
2 Randomly initialize encoder/decoder params.  $\theta_{\text{enc.}}^0, \theta_{\text{dec.}}^0$ 
3 Clamp pre-trained model  $f(\cdot; \theta_{\text{task}})$  params.  $\theta_{\text{task}}$ 
4 for  $\tau \leftarrow 0$  to  $T$  do
5   Sample  $\{(x_i, y_i)\}_{i=1}^B$  from dataset  $\mathbf{D}$  with
      $y = f(x; \theta_{\text{task}})$ 
6   Encode  $z = p(z | x; \theta_{\text{enc.}}^\tau)$ 
7   Decode  $\hat{x} = q(x | z; \theta_{\text{dec.}}^\tau)$ 
8   Compute Predictions  $\hat{y} = f(\hat{x}; \theta_{\text{task}})$ 
9    $\theta_{\text{enc.}}^{\tau+1}, \theta_{\text{dec.}}^{\tau+1} \leftarrow$ 
     BACKPROP[ $\mathcal{L}_{\text{weight}}(x, \hat{x}, y, \hat{y}; \theta_{\text{task}}, \theta_{\text{enc.}}^\tau, \theta_{\text{dec.}}^\tau)$ ]
10 end
Result: Return learned parameters  $\theta_{\text{enc.}}^T, \theta_{\text{dec.}}^T$ .

```

---

#### 3.1 Wide applicability of our Co-design Algorithm

While simple, Alg. 1 is a powerful, general solution for task-relevant communication. A principal benefit is that we can either invoke a pre-trained task module fully remotely, or also flexibly adjust the split of computation between a resource-constrained robot and server. In the latter scenario, depicted in Fig. 5a, we can take a large, pre-trained DNN and only execute a fraction of layers locally on the robot to flexibly adjust for resource constraints. Then, we can insert a minimal, learned encoder  $p$  to compress an intermediate result, transmit bottleneck  $z$ , and then decode with  $q$  to continue computation with the bulk of the pre-trained model at a server.

Importantly, *only* the distributed encoder and decoder are trained, unlike the fixed task module. Further, by adjusting reconstruction weight  $\lambda$ , which serves as a regularizer, we can significantly reduce reconstruction error with only a marginal gain in bottleneck size  $z$ . Our ability to flexibly set  $\lambda$  leads to the following variants of Alg. 1:

1. Task-Aware-Weighted (TASKNET): Our core contribution is the weighted task-aware training scheme, henceforth referred to as TASKNET. By accounting for the pre-trained task network  $f(\cdot; \theta_{\text{task}})$  when encoding sensory inputs, TASKNET achieves low task loss for much less transmitted data  $z$ . We emphasize that TASKNET refers to our co-design training scheme and is compatible with any domain-specific, pre-trained task module  $f$ . TASKNET’s reconstruction loss weight  $\lambda$  (Eq. 1) can be flexibly set by a roboticist per application. All our experiments performed a scan over evenly-spaced values of  $\lambda$  and only plotted one representative curve, in orange, for visual clarity.
2. Fully Task-Aware: To quantify the full compression benefits of optimizing only for an end-task, we consider the special case of TASKNET when  $\lambda = 0$ , which is suitable for scenarios when video inputs are automatically classified in real-time by a DNN and not intended for a human viewer. This scheme, always colored in blue, achieves the lowest task loss for the smallest representation size  $Z$ , which is ideal for scenarios with strict network bandwidth limits. We now evaluate both variants of Alg. 1, as well as benchmark schemes, on large-scale deep learning experiments.

## 4 Experimental Results

To evaluate TASKNET, we first introduce two benchmark algorithms and then describe common evaluation metrics for all schemes. Then, we highlight the wide utility of TASKNET across various sensing modalities and task DNN modules, for tasks ranging from Martian terrain classification to motion planning and environmental timeseries classification. Further, we benchmark performance on the standard MNIST dataset. To recreate our work, we provide all software, data, and pre-trained DNNs for the Edge Tensor Processing Unit (TPU) and servers at the following website:

<https://sites.google.com/view/task-net>.

### 4.1 Evaluation Metrics and Benchmark Algorithms

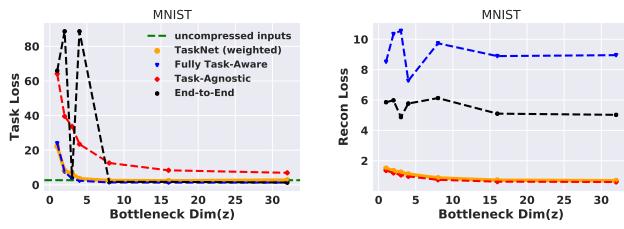
The principal objective of Prob. 1 is to minimize task loss, and optionally minimize weighted reconstruction loss, in scenarios where decoded inputs  $\hat{x}$  need to be inspected. As such, we show that TASKNET achieves (A) low task loss (i.e. high task accuracy) and (B) flexibly achieves low reconstruction loss for small bottleneck sizes  $z$ . To assess the effects of compression on task loss, we quantify the task loss achieved by

passing uncompressed, *original* sensory inputs  $x$ , without an information bottleneck, into the task network  $f$ . This metric, henceforth referred to as the *uncompressed input task loss*, represents the best task accuracy we can achieve without network constraints, and is pictured by the green dashed line in all task loss figures (e.g. Fig. 3). Further, all figures show results on each domain’s *test* dataset. We test the above metrics on the following two benchmarks that represent conventional paradigms for networked perception:

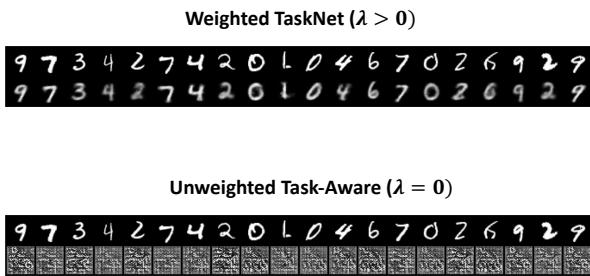
1. Task-Agnostic: For each bottleneck dimension  $Z$ , this scheme trains a standard variational autoencoder to minimize reconstruction loss, and simply passes the decoded input  $\hat{x}$  through the pre-trained task module  $f(\cdot; \theta_{\text{task}})$ . The task-agnostic scheme is always colored in red in all figures.
2. End-to-End: This scheme has the exact same encoder, decoder, and task module architectures as the others, but allows the task parameters  $\theta_{\text{task}}$  to be updated. In principle, an end-to-end approach can achieve the same compression and accuracy as our task-aware co-design approach since all parameters can be optimized for the ultimate weighted task loss. However, in practice, all end-to-end schemes took longer to train and performed worse than TASKNET for small bottlenecks  $z$  when the networks got stuck in local minima during training, such as in Fig. 3 for  $Z = 2$  and  $Z = 4$ . In essence, an end-to-end approach does not exploit the structure of task network  $f$  for faster, more stable representation learning. We now evaluate all schemes.

### 4.2 Digit Classification with an Information Bottleneck

We first benchmark TASKNET on the standard MNIST [24] dataset, where task network  $f(\cdot; \theta_{\text{task}})$  was a publicly-available ResNet DNN for digit classification [11]. Fig. 3a illustrates that TASKNET (orange) achieves the same lower-bound classification error as passing *original, unperturbed* images to the task DNN, even for images that are compressed with an extremely small bottleneck of size  $Z = 3$ . In contrast, the task-agnostic scheme (red) is still not able to achieve the lower bound classification error even for a much larger bottleneck of  $Z = 32$ , showing that TASKNET improves the compression ratio by  $10\times$  over a standard autoencoder with *better* task accuracy. Fig. 3b shows that TASKNET decodes images  $\hat{x}$  with low reconstruction loss, which yields human-interpretable, but highly compressible, decoded samples that are visualized in Fig. 4. While we did a sweep over reconstruction weights  $\lambda$ , we plot results for  $\lambda = 0.01$ , which was chosen to reflect similar emphasis on compressed accuracy and reconstruction loss. Further, our learned encoder  $p$  and decoder  $q$  were standard 3-layer convolutional variational autoencoders that added a small overhead of only  $< 7\%$  extra compute compared to solely running task network  $f$ . In the end-to-end scheme, the task network is also learned during training.



**Fig. 3: MNIST test dataset results:** TASKNET achieves the same lower-bound classification error as passing *original* images through the task network (green) for a small bottleneck  $Z = 3$ . In contrast, a task-agnostic method, in red, obtains worse classification error for even larger  $Z = 32$  since it does not emphasize salient features for classification.



**Fig. 4: Motivation for regularization with a reconstruction loss.** The key benefit of adding a reconstruction weight with  $\lambda > 0$  in our task loss (Eq. 1) is to yield decoded images that are human-interpretable, as shown for our warm-up MNIST example. In each panel, we show original images  $x$  in the top row and decoded images  $\hat{x}$  in the bottom row. **Top Panel, TASKNET with  $\lambda = 0.01$ :** Clearly, adding a regularization term with a representative reconstruction weight  $\lambda = 0.01$  in Eq. 1 learns representations that are highly-compressible with  $Z = 4$ , but yield interpretable decoded samples  $\hat{x}$  that a human can debug. **Bottom Panel,  $\lambda = 0$ :** In contrast, if we purely optimize for a machine’s task loss with reconstruction weight  $\lambda = 0.0$ , decoded images  $\hat{x}$  (bottom row) are *not* human-interpretable, motivating our weighted co-design approach.

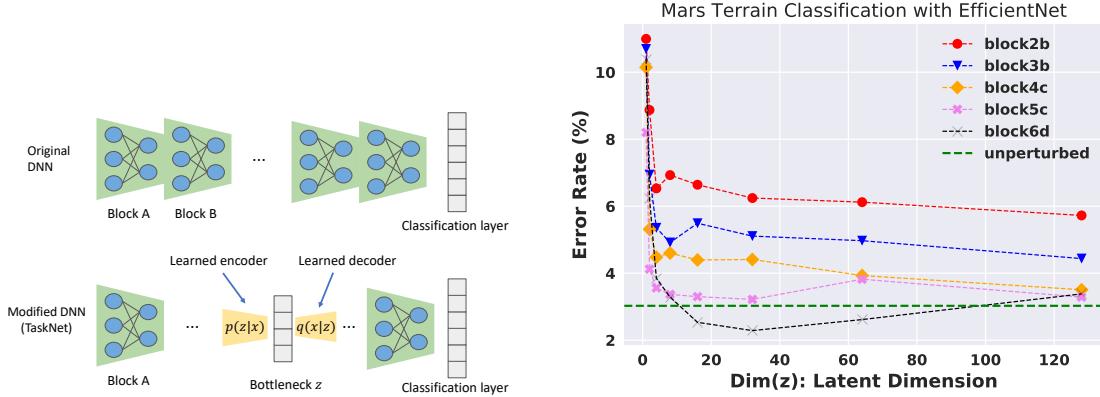
#### 4.3 Compute-Efficient Mars Terrain Classification

To test TASKNET with complex vision DNNs, we consider a scenario where a low-power Mars rover encodes terrain images for remote classification at a powerful orbiting compute server, a scenario inspired by recent proposals [13]. We used Martian terrain images from NASA’s HiRISE Dataset [9], consisting of 8 classes including dunes, craters, impact ejecta, and other formations. The labeled HiRISE Dataset images [9] are several hundreds of megabytes large, unlike the larger HiRISE panoramas we use for our subsequent motion planning experiments in Sec. 4.4. Task network  $f(\cdot; \theta_{\text{task}})$  is an EfficientNet-B0 DNN, which is the smallest of a family of DNNs that trade off accuracy and model complexity [36]. Our EfficientNet has 5.36 million parameters and classification accuracy of 95.7% on 7303 HiRISE test images.

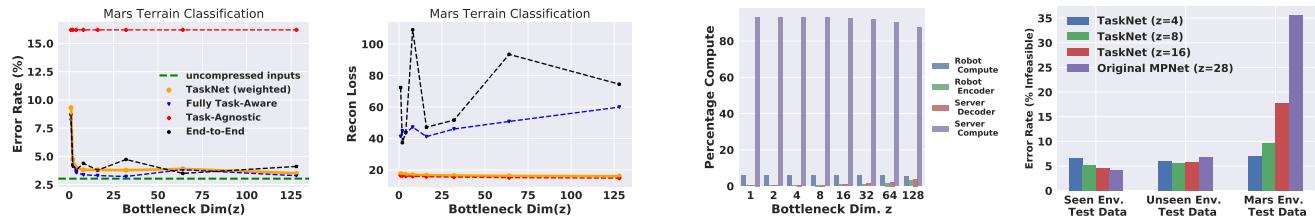
Given the relatively large size of the vision DNN, we want to flexibly split computation between a compute-limited robot and server. Notably, as shown in Fig. 5a, we can exploit a common architecture in many modern DNNs, which are composed of *sequential blocks* of several convolutional layers, to adjust the extent of on-robot compute. For example, Fig. 5a shows how EfficientNet block-4c consists of all parameters from the input to a specific layer  $l$  and block-5c has all successive parameters until layer  $l' > l$ . We first train the full, modular EfficientNet, which we use as our pre-trained model,  $f(x; \theta_{\text{task}})$ , as mentioned in Alg. 1. Then, we constrain the robot to only run part of the large DNN until a specific block  $b$ , which yields an intermediate convolutional feature map, and then insert our minimal task encoder  $p(\cdot; \theta_{\text{enc}})$  to compress the feature map to encoding  $z$ . Encoding  $z$  is transmitted to a remote server, which uses our task decoder  $q(\cdot; \theta_{\text{dec}})$  to reconstruct the original feature map generated by block  $b$  and resumes computation for a final classification. Importantly, only the encoder/decoder are learned with the *fixed* EfficientNet. For this experiment, the value of  $\lambda$  doesn’t matter as the primary focus is not on reconstruction, as the reconstruction of the hidden layer output is not interpretable. We set  $\lambda = 0$ .

Fig. 5b details several configurations where the robot only encodes inputs up to EfficientNet block  $b$ , but still achieves within 1.5% of the uncompressed input task loss. For example, Figs. 5b and 6 show by computing until block-4c, TASKNET achieves  $3.87 \times$  lower classification error than conventional task-agnostic methods, with a small bottleneck of  $Z = 128$  32-bit floating-point values. This enables our scheme to send only 3.27% of data compared to the full feature map of size  $14 \times 14 \times 80$ . However, by introducing our learned robot encoder  $p$  and server decoder  $q$  (pictured in Fig. 5a), we introduce a small overhead of  $< 7.37\%$  extra compute compared to simply running the original EfficientNet, as shown in Fig. 7a. In practice, this small compute overhead is tolerable given the accuracy benefits of querying a powerful remote server with a small transmitted data representation. Additionally, in Fig. 5b we observe that when mapping the intermediate feature map from block 6d, for lower bottleneck dimensions  $Z$ , the task error rate is lower than the unperturbed case. This is an artifact of training DNNs and there is no specific reason for this trend.

**Robot compute savings:** Notably, by only devoting computational effort to encode task-relevant information, TASKNET reduces robot computation to only 8.5% of overall compute, compared to the 91.5 % delegated to a remote server (Fig. 7a). As such, TASKNET is significantly more efficient than naively just splitting a compute-intensive DNN in half or running it fully on-board, which could be infeasible for micro-robots or miniature satellites like KickSats [27].



**Fig. 5: Mars terrain classification. (Left) Flexible allocations of robot and server compute:** In a nominal configuration, our co-design algorithm fully runs a modular perception model at a remote server and simply encodes task-relevant features at a robot, as shown in Fig. 1. However, a key benefit of our approach is that we can *also* flexibly divide computation for a large perception model between a robot and server. As shown above, we can take a *pre-trained* EfficientNet and insert a learned, task-relevant encoder (yellow) to map an intermediate feature map to a bottleneck representation  $z$ . Once  $z$  is transmitted, the learned decoder (yellow) uncompresses the intermediate feature map so that computation can continue at a remote server with the rest of the pre-trained EfficientNet. Unlike classical approaches that simply split large DNNs between devices [16, 15], we co-learn an encoder/decoder (yellow) along with the pre-trained parameters (green) for high accuracy with a small  $Z$ . **(Right)** We plot the accuracy for various pictured allocations between robot and server compute. Specifically, we plot TASKNET with  $\lambda = 0$  for scenarios where the robot only runs up to EfficientNet block  $b$  locally and achieve within 1.5% of the uncompressed input error rate (green dashed line) at block 4c and beyond for  $Z = 128$ .



**Fig. 6: Terrain classification.** TASKNET achieves much lower classification error than task-agnostic methods (red), where even a bottleneck dimension of  $z = 128$  corresponds to sending only 3.27 % of data compared to original convolutional feature maps. While our method incurs a compute overhead by adding a task-relevant encoder and decoder, it is less than 7.37% of total computation. Such overhead is worthwhile for robots to unlock the accuracy of remote computation while transmitting minimal data.

#### 4.4 Robotic Motion Planning with Neural Networks

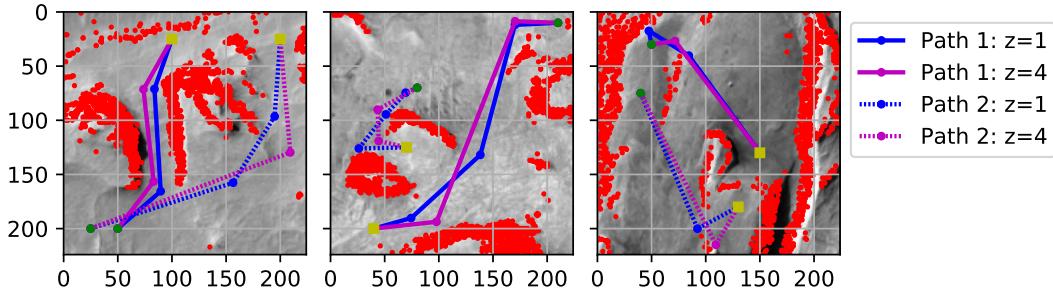
To demonstrate TASKNET’s utility for motion planning, we consider a scenario where a low-power drone takes an aerial terrain image and sends a compressed representation to a remote server, which generates a collision-free motion plan for a swarm of ground robots. Our task network  $f(\cdot; \theta_{\text{task}})$  is a *Motion Planning Network* (MPNET), a recently published planner [33] that generates collision-free paths an order of magnitude faster than conventional sampling-based planners such as RRT\* [17], but with virtually the same path cost. MPNET consists of an obstacle encoder network, such as

(a) Martian Terrain Classification      (b) Neural Motion Planning

**Fig. 7:** (a) By only encoding task-relevant information, TASKNET reduces on-robot computation compared to what is delegated to a remote server (purple). (b) TASKNET generalizes better than the original MPNET on out-of-domain Mars terrain data to find more feasible paths (y-axis).

a contractive autoencoder, which maps an obstacle point-cloud to an embedding of dimension up to  $Z = 28$ . Then, MPNET’s planning network maps the obstacle embedding  $z$  as well as current and goal robot configurations to predict the next navigation waypoint. The original MPNET DNN planners for two and three dimensional environments are up to 41 MB and 69 MB in size. While of moderate size compared to large perception DNNs, they are too large to run on low-power accelerators like the Edge TPU, which runs only compressed DNNs that fit in 8MB of on-chip memory.

In our experiments, the MPNET encoder runs on the robot, but the *pre-trained* planning network runs on a remote server to potentially plan paths for a robotic swarm. Crucially,



**Fig. 8: Neural Motion Planning on Mars Terrain Data:** Our task-aware motion planner generalizes to unseen Martian terrain data for much smaller point-cloud embedding sizes  $z$  compared to a publicly-available neural motion planner.

rather than transmitting a large obstacle embedding  $z_0$  as in the original MPNET, we insert a small encoder at the robot which generates a minimal representation  $z \ll z_0$ , which is sent across a wireless link, decoded to  $z_0$ , and passed into the planning network to generate a motion plan. As described in Alg. 1, we only learn the encoder and decoder to generate compressed environment encoding  $z$  and use the original MPNET planning network. Since the original embedding  $z_0$  is not human-interpretable, we plot results for  $\lambda = 0$  though we thoroughly evaluated  $\lambda$  uniformly. Our task-relevant encoder/decoder consist of a small fully connected neural network, leading to a compute overhead of only  $< 2\%$  more parameters.

Fig. 7b shows the percentage of collision-free, feasible paths (y-axis) generated by our TASKNET co-designed planner and the original MPNET on the same 110 test environments and 40K test paths provided by the MPNET authors. In all cases, if a feasible, collision-free path is found, its path length is within 1% of the path generated by the RRT\* planner. Our key result is that TASKNET achieves collision-free paths with costs and failure rate within 1% of the original MPNET, but with much smaller environment encodings  $z$ . For example, TASKNET achieves good performance for  $z = 4$  as opposed to  $z = 28$  for the original MPNET, leading to a compression gain of  $7\times$ .

#### Generalization of motion planner to unseen domains:

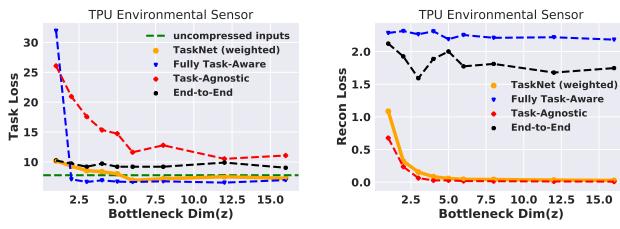
We further tested the TASKNET motion planner on obstacle point clouds from the Mars HiRISE dataset, where the red point clouds in Fig. 8 represent regions above 30 degrees of elevation. Even though MPNET and TASKNET's train and test environments had point clouds representing *polygonal* obstacles, TASKNET adapts to curved contours of *never-before-seen* Mars terrain with small bottlenecks  $z$ .

Fig. 7b validates the above observation on over 40.5K paths, where both MPNET and TASKNET performed similarly on the first two scenarios of seen and unseen environments in MPNET's original test data. Specifically, a seen test environment was observed by the planner during training, albeit with different start and end test configurations [33].

Our key experimental result is shown in the third column of Fig. 7b, where TASKNET significantly outperforms the original MPNET for small bottlenecks  $z = 4, 8$  on out-of-domain Mars terrain data. We hypothesize that TASKNET generalizes better for small  $z$  since it focuses the limited information capacity of the bottleneck to represent salient features for planning, and sacrifices reconstruction loss, which is irrelevant to the task goal. Since we only had access to a large, out-of-domain dataset for the Mars example, we plan to further stress test TASKNET's generalization capabilities in future work.

#### 4.5 Environmental sensor timeseries anomaly detection

To highlight how TASKNET generalizes to multiple sensor modalities, we emulated the scenario of a micro-robot, potentially part of a swarm, that streams environmental sensory data to a remote server for anomaly detection. We used an environmental sensor that connects to the Edge TPU to measure light, temperature, pressure, and humidity timeseries. The task network  $f(\cdot; \theta_{\text{task}})$  is a compact neural network that classifies each timeseries window of length  $w$ ,  $x^{t-w:t}$ , into three sensor conditions  $y^t$ . The first class  $y^t = 0$  represents if the sensor is being tightly clamped and tampered with, which leads to a rapid fluctuation in the light and humidity measurements. The second class  $y^t = 1$  represents when the sensor reads a temperature spike, for which we generate training data by placing a hot hair-dryer over the sensor. The third class  $y^t = 2$  represents natural environmental variation without anomalies. We collected two weeks of sensor data for a total of 30 training and 30 test traces of 5 minutes each, equally balanced across all three classes. The task network  $f(\cdot; \theta_{\text{task}})$ , which was trained on stochastic, diverse timeseries, achieved a promising 90% test accuracy. Fig. 9 shows that TASKNET aggressively compresses the timeseries measurements and outperforms the task-agnostic benchmark. Furthermore, we plot for  $\lambda = 1$  which reflects similar weight between reconstruction loss and accuracy. Overall, our diverse results show



**Fig. 9: Environmental Sensing:** TASKNET achieves low reconstruction loss and matches the low classification error rate achieved when using *original* sensor data.

the promise of provisioning publicly-available task modules at a central server and co-designing minimal representations for networked perception.

#### 4.6 Image Classification on CIFAR-10 with Information Bottleneck

In this section, we focus on benchmarking the TASKNET performance on real-world images from the CIFAR-10 dataset [22]. The task network  $f(\cdot; \theta_{\text{task}})$  is a ResNet-50 DNN which provides an upper-bound classification accuracy of 93.65%. We use a standard 3-layer Convolutional Autoencoder for our encoder and decoder. Fig. 10 illustrates the performance gain of TASKNET in comparison with Task-Agnostic co-design, in terms of lower task loss and higher task accuracy for different bottleneck ( $Z$ ) dimensions on the CIFAR-10 test dataset.

The value of the reconstruction weight,  $\lambda$ , had to be chosen separately for each  $Z$ . This choice of  $\lambda$  can be inferred directly from the test accuracy plot in Fig. 10, as follows. For  $Z \leq 256$ , we set  $\lambda$  to 0.01, as the test accuracy depends solely on minimizing the task loss, and the reconstruction loss is only a small part to ensure meaningful/human-interpretable decoded images  $\hat{x}$ . Whereas, for  $Z > 256$ , we set  $\lambda$  as 100, as the test accuracy mainly depends on accurate reconstruction.

We also note that the task loss and thereby task accuracy remain fairly constant for varying bottleneck dimensions, for TASKNET and Fully Task-Aware co-design. The same trend is also observed in the MNIST benchmark. This further strengthens our case of task-relevant representation. This trend can be explained intuitively as follows: for 10 target classes (as is the case in MNIST and CIFAR-10), it is sufficient to have a small encoding space to learn class-specific features such that the task performance is maximized.

Fig. 11 shows the decoded images for TASKNET and Task-Agnostic co-design for  $Z = 256$ .  $Z = 256$  provides 82.89% classification accuracy on the test set (accuracy drop of only around 10% from the upper bound) along with human-interpretable decoded images. The additional compute overhead for our autoencoder architecture is only 14.5% when compared to solely running the task network. Given that the

input image has 3072 8-bit unsigned integer values, transferring only 256 32-bit floating-point values provides 66.67% compression gain. For the same range of task accuracy, the compression gain can be as high as 97.9% for  $Z = 16$ , if human-interpretability is not required. In comparison with  $Z = 256$  for Task-Agnostic approach, both TASKNET and Fully Task-Aware co-design provided test dataset accuracies greater than 80% even for  $Z = 16$  leading to  $16\times$  relative compression gain.

#### 4.7 Human-interpretability with adversarially trained task network co-design

Though the Fully Task-Aware co-design ( $\lambda = 0$ ) learns task-relevant features, from Fig. 4 and Fig. 11, we see that the decoded image  $\hat{x}$  is not human-interpretable. The task network is trained to minimize a certain task objective function for a given dataset, and hence, there is no constraint on the features learned by the task network to be human-aligned or human-interpretable. Extending this notion to our Fully Task-Aware co-design, the Autoencoder parameters are optimized to map any image  $x$  to a decoded image  $\hat{x}$  that is not human-interpretable, but maximizes the task performance. However, since the compression gains are high, we explore the possibility of obtaining human-interpretable decoded images for the Fully Task-Aware co-design, i.e. human-interpretability without reconstruction loss.

In particular, we focus on adversarial training of the task network as explained in [39], which suggests that one of the advantages of training the task network adversarially is that the learned features are human-perceptible. Adversarial training minimizes the expected adversarial loss as specified in [39],

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \max_{\delta \in \Delta} L(x + \delta, y; \theta_{\text{task}}) \right], \quad (5)$$

where  $x$  and  $y$  are the input and target from the dataset  $\mathcal{D}$ ,  $\Delta$  is the set of  $l_2$  bounded perturbations and  $\theta_{\text{task}}$  represents the parameters of the adversarially trained task network  $f_{\text{adv}}(\cdot; \theta_{\text{task}})$ .

We hypothesize that if the task network is adversarially trained, then even for the Fully Task-Aware co-design, the decoded image  $\hat{x}$  should be human-interpretable. Instead of minimizing the task loss, we aim to minimize the distance between the input image's and the decoded image's *representation*. As in [12], we define the *representation*  $R(x)$  as the activations of the penultimate layer of the task network. In the case of ResNet-50 for CIFAR-10,  $R(x) \in \mathbb{R}^{512}$ . What follows the penultimate layer is a linear layer that takes the representation as input to generate logits as output. Therefore, minimizing the *representation* distance is the same as minimizing the task loss. In fact, we note that it is better

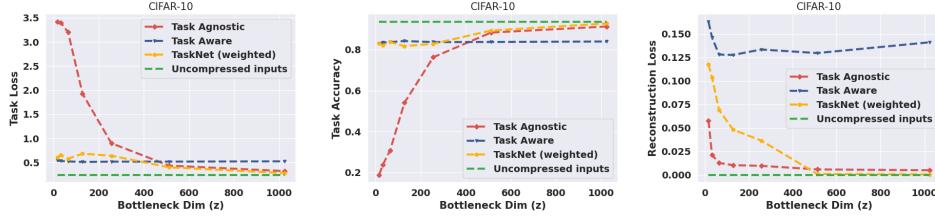


Fig. 10: **CIFAR-10 test dataset results.** The plot on the left and middle show the variation of classification loss and classification accuracy respectively, on the test dataset for different  $Z$  values. TASKNET (yellow) achieves lower classification loss and higher classification accuracy for all  $Z$  values with respect to the Task-Agnostic co-design (red). TASKNET is able to achieve 82.89% accuracy for  $Z = 256$ , which is only 10% lower than the upper-bound accuracy. The plot on the right shows how the reconstruction loss varies with  $Z$  values. Here again, TASKNET achieves similar or even lower reconstruction loss as in the Task-Aware case, given sufficient  $Z$  size.

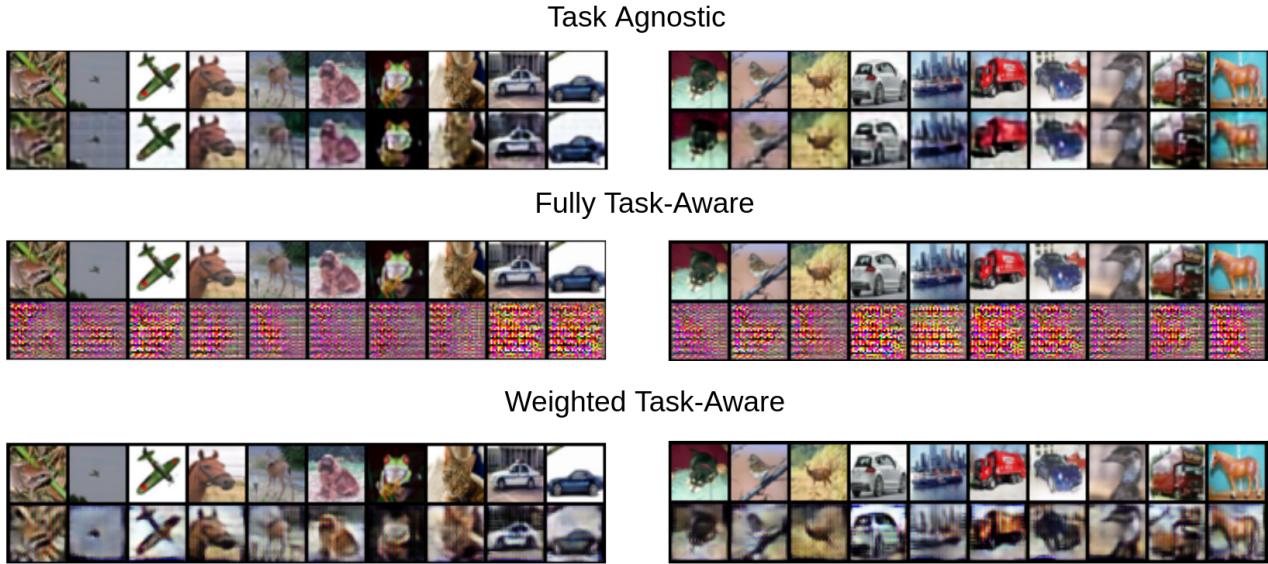
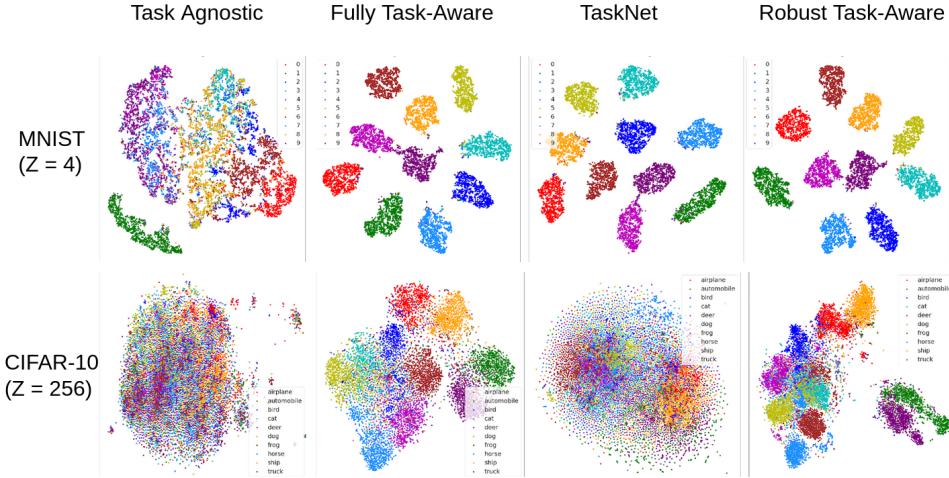


Fig. 11: **CIFAR-10 classification examples.** Shown here are some qualitative results of Task-Agnostic, Fully Task-Aware and TASKNET (Weighted Task-Aware,  $\lambda = 0.01$ ) co-designs on the CIFAR-10 test dataset for  $Z = 256$ . For each co-design, the first row is the input image and the second row is the decoded output. It can be noted that the decoded outputs are human-interpretable for TASKNET, except for color shifts with respect to the input image. The color shifts are inevitable consequences of arbitrary color-jitter data augmentations used to improve the TASKNET generalization abilities.

to minimize the *representation* distance in terms of obtaining human-interpretable decoded images. This is because, while updating the encoder and decoder parameters to minimize the task loss, the loss gradients have to be backpropagated through linear layers in the task network, which destroy the spatial structure. Fig. 13 shows the stark increase in the decoded image’s interpretability for Fully Task-Aware co-design using an adversarially trained task network in comparison with Fig. 11 (row 2, Fully Task-Aware). It also shows the subtle difference between minimizing the task loss and the *representation* distance, in terms of human-interpretability.

	Task Agnostic	Fully Task-Aware	TASKNET (weighted)	Robust Task-Aware
MNIST	0.05	<b>0.57</b>	0.50	<b>0.57</b>
$Z = 4$				
CIFAR-10	-0.06	0.07	-0.04	<b>0.10</b>
$Z = 256$				

Table 1: **Silhouette Coefficient.** Clustering scores of the embedding space for different methods clearly show that the task-aware co-designs perform much better. In particular, the robust task-aware co-design stands out, as evident in Fig. 12 as well. Robust Task-Aware here refers to the Fully Task-Aware co-design with adversarially trained task network. The score is bounded between -1 for incorrect clustering and +1 for highly dense clustering respectively. Scores around zero indicate overlapping clusters.



**Fig. 12: Clustering of Embedding space  $Z$ .** The introduction of task loss during the Autoencoder training leads to significant improvement in the clustering of data points in the embedding space. There is a progressive increase in class separation as the loss changes from Task-Agnostic to TASKNET and then to Fully Task-Aware co-designs. This is illustrated for both the MNIST and CIFAR-10 datasets. Robust Task-Aware refers to the Fully Task-Aware co-design with adversarially trained task network. Robust Task-Aware offers the best separation between classes.



**Fig. 13: Fully Task-Aware co-design with adversarially trained task network:** The figure shows the decoded outputs of Fully Task-Aware co-design using an adversarially trained task network. The results for minimizing the task loss and minimizing the *representation* distance are shown. In both cases, the first is the input and the second row is the decoded output. Minimizing the *representation* distance provides human-interpretable decoded output. For example, we can observe the face of the horse, the texture and pose of frog, horns of the deer, the shapes and backgrounds of airplanes and ships, etc. in the decoded output.

#### 4.8 Improved clustering in embedding space with task-aware co-design

One of the benefits of optimizing the autoencoder parameters for task performance is a bottleneck space with data points clustered according to the class label. In other words, inputs of the same target class group together in dense clusters and those from different classes are relatively well separated. This can be observed in Fig. 12 which shows  $Z$  reduced to two dimensions using t-SNE visualization [41]. The Fully Task-Aware co-design results in better clustering of inputs in the embedding space for both MNIST and CIFAR-10 datasets. To quantify the clustering performance of the encoders we utilize Silhouette Coefficients [34]. The silhouette coefficient is a measure of how close a sample, on average,

is to its own cluster compared to other clusters. The values in Table. 1 indeed corroborate the qualitative results from t-SNE visualization. Between the two Fully Task-Aware co-designs, the clustering with adversarially trained task network  $f_{adv}(\cdot; \theta_{task})$  is better overall. Given the advantageous properties of co-designing an autoencoder with an adversarially pre-trained task network, further investigation is merited.

*Limitations of our work:* Our current method requires a constraint on the maximum bottleneck  $Z$ , based on wireless network capacity, and performs a search over increasing sizes  $z$ . A promising future direction is to automatically determine the optimal representation dimension  $Z^*$  that minimizes weighted loss. Additionally, our encoder and decoder parameters need to be updated if the task network is modified. This occurs very frequently in real-world applications

where the deployed DNNs are constantly updated with recent/new training samples. The effect of communication loss is another aspect that was not the focus of this work. It will be interesting to observe the effects on task performance, in the presence of communication loss, when an adversarially trained task network is used. Further, we could extend our experiments to compress LIDAR point clouds and long segments of video.

## 5 Discussion and Conclusions

This paper presents a novel framework to aggressively compress rich robotic sensory data for the ultimate needs of a *machine* sensing task, which allows robots to reduce on-board computation and communication bandwidth by up to  $11 \times$ . A key benefit of our co-design algorithm is that it is rooted in linear systems theory and gracefully scales to complex DNNs. We envision that our co-design algorithm can be applied to a variety of resource-constrained robots that need to distill sensory data for remote inference, such as future Mars Rovers, microrobots, and even low-power drones that collaborate by communicating over 5G wireless networks. Additionally, we also show that the compressed form of the rich sensory data can be human-interpretable if the task network is adversarially trained.

In future work, we plan to extend our algorithm to create minimal, task-relevant representations for cooperative control and develop representations that generalize across several tasks, leveraging ideas from multi-task and meta-learning. While our approach reduces on-board computation, we plan to also quantify the power consumption of communicating minimal representations. Finally, given our promising results with EfficientNets and deep learning accelerators, we plan to deploy our algorithm on a networked robotic autonomy stack which supports remote inference and tele-operation. We also plan to experiment with communication loss and its interplay with robust/adversarially trained task networks.

## 6 Declarations

### 6.1 Funding

The NASA University Leadership initiative (grant #80NSS-C20M0163) and Lockheed Martin Corporation provided funds to assist the authors with their research, but this article solely reflects the opinions and conclusions of its authors and not any NASA or Lockheed Martin Corp. entity.

### 6.2 Competing Interests

The authors have no competing interests to declare that are relevant to the content of this article.

## References

1. Edge tpu. <https://cloud.google.com/edge-tpu/>, 2019. [Online; accessed 01-Sep.-2019].
2. Mars reconnaissance orbiter: Communications with earth. <https://mars.nasa.gov/mro/mission/communications/>, 2020. [Online; accessed 18-Oct.-2020].
3. Semantic distortion networks. <https://on-demand.gputechconf.com/gtc/cdc/2018/pdf/dc8211-semantic-distortion-networks-and-embedding-large-scale-semantic-graphs-presented-by-lockheed-martin.pdf>, 2022.
4. Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. A survey of model compression and acceleration for deep neural networks. *arXiv preprint arXiv:1710.09282*, 2017.
5. Sandeep Chinchali\*, Eyal Cidon\*, Evgenya Pergament\*, Tianshu Chu, and Sachin Katti. Neural networks meet physical networks: Distributed inference between edge devices and the cloud. In *Proceedings of the 17th ACM Workshop on Hot Topics in Networks*, HotNets '18, pages 50–56, 2018. \* Indicates equal contribution.
6. Sandeep Chinchali, Apoorva Sharma, James Harrison, Amine Elhafsi, Daniel Kang, Evgenya Pergament, Eyal Cidon, Sachin Katti, and Marco Pavone. Network offloading policies for cloud robotics: a learning-based approach. *arXiv preprint arXiv:1902.05703*, 2019.
7. Michael Crawshaw. Multi-task learning with deep neural networks: A survey. *arXiv preprint arXiv:2009.09796*, 2020.
8. Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
9. Gary Doran, Steven Lu, Lukas Mandrake, and Kiri Wagstaff. Mars orbital image (HiRISE) labeled data set version 3, January 2019. URL <https://doi.org/10.5281/zenodo.2538136>.
10. John Emmons, Sadjad Fouladi, Ganesh Ananthanarayanan, Shivaram Venkataraman, Silvio Savarese, and Keith Winstein. Cracking open the dnn blackbox: Video analytics with dnns across the camera-cloud boundary. In *Proceedings of the 2019 Workshop on Hot Topics in Video Analytics and Intelligent Edges*, pages 27–32, 2019.
11. Logan Engstrom, Andrew Ilyas, Hadi Salman, Shibani Santurkar, and Dimitris Tsipras. Robustness (python library), 2019. URL <https://github.com/MadryLab/robustness>.
12. Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Brandon Tran, and Aleksander Madry. Adversarial robustness as a prior for learned representations.

- arXiv preprint arXiv:1906.00945*, 2019.
13. Keck Institute for Space Studies. Virtual workshop: Nebulae - deep-space computing clouds - part 2, 2020. URL <https://www.kiss.caltech.edu/workshops/nebulae/nebulae2.html#Schedule>.
  14. Muhammad Ghifary, W Bastiaan Kleijn, Mengjie Zhang, and David Balduzzi. Domain generalization for object recognition with multi-task autoencoders. In *Proceedings of the IEEE international conference on computer vision*, pages 2551–2559, 2015.
  15. Babak Hassibi, David G Stork, and Gregory J Wolff. Optimal brain surgeon and general network pruning. In *IEEE international conference on neural networks*, pages 293–299. IEEE, 1993.
  16. Yiping Kang, Johann Hauswald, Cao Gao, Austin Rovinski, Trevor Mudge, Jason Mars, and Lingjia Tang. Neurorurgeon: Collaborative intelligence between the cloud and mobile edge. *ACM SIGPLAN Notices*, 52(4):615–629, 2017.
  17. Sertac Karaman, Matthew R Walter, Alejandro Perez, Emilio Frazzoli, and Seth Teller. Anytime motion planning using the rrt. In *2011 IEEE International Conference on Robotics and Automation*, pages 1478–1483. IEEE, 2011.
  18. Abdallah Kassir, Robert Fitch, and Salah Sukkarieh. Communication-aware information gathering with dynamic information flow. *The International Journal of Robotics Research*, 34(2):173–200, 2015.
  19. Ben Kehoe, Akihiro Matsukawa, Sal Candido, James Kuffner, and Ken Goldberg. Cloud-based robot grasping with the google object recognition engine. In *Robotics and Automation (ICRA), 2013 IEEE International Conference on*, pages 4263–4270. IEEE, 2013.
  20. Ben Kehoe, Sachin Patil, Pieter Abbeel, and Ken Goldberg. A survey of research on cloud robotics and automation. *IEEE Trans. Automation Science and Engineering*, 12(2):398–409, 2015.
  21. Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
  22. Alex Krizhevsky. Learning multiple layers of features from tiny images. pages 32–33, 2009. URL <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>.
  23. J Kuffner. Cloud-enabled robots in: Ieee-ras international conference on humanoid robots. *Piscataway, NJ: IEEE*, 2010.
  24. Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>.
  25. Pusong Li, Bill DeRose, Jeffrey Mahler, Juan Aparicio Ojea, Ajay Kumar Tanwani, and Ken Goldberg. Dex-net as a service (dnaas): A cloud-based robust robot grasp planning system. In *2018 IEEE 14th International Conference on Automation Science and Engineering (CASE)*, pages 1420–1427. IEEE, 2018.
  26. Zihao Liu, Tao Liu, Wujie Wen, Lei Jiang, Jie Xu, Yanzhi Wang, and Gang Quan. Deepn-jpeg: A deep neural network favorable jpeg-based image compression framework. In *Proceedings of the 55th Annual Design Automation Conference*, pages 1–6, 2018.
  27. Zachary Manchester, Mason Peck, and Andrew Filo. KickSat: A crowd-funded mission to demonstrate the world’s smallest spacecraft. 2013.
  28. Gajamohan Mohanrajah, Vladyslav Usenko, Mayank Singh, Raffaello D’Andrea, and Markus Waibel. Cloud-based collaborative 3d mapping in real-time with low-cost robots. *IEEE Transactions on Automation Science and Engineering*, 2015.
  29. Manabu Nakanoya, Sandeep Chinchali, Alexandros Anemogiannis, Akul Datta, Sachin Katti, and Marco Pavone. Co-design of communication and machine inference for cloud robotics. *Robotics: Science and Systems XVII, Virtual Event*, 2021.
  30. Fabrizio Nenci, Luciano Spinello, and Cyrill Stachniss. Effective compression of range data streams for remote robot operations using h. 264. In *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3794–3799. IEEE, 2014.
  31. Vincent Pacelli and Anirudha Majumdar. Task-driven estimation and control via information bottlenecks. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 2061–2067. IEEE, 2019.
  32. Vincent Pacelli and Anirudha Majumdar. Learning task-driven control policies via information bottlenecks. *arXiv preprint arXiv:2002.01428*, 2020.
  33. Ahmed H Qureshi, Anthony Simeonov, Mayur J Bency, and Michael C Yip. Motion planning networks. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 2118–2124. IEEE, 2019.
  34. Peter J. Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20:53–65, 1987.
  35. Anoopkumar Sonar, Vincent Pacelli, and Anirudha Majumdar. Invariant policy optimization: Towards stronger generalization in reinforcement learning. *arXiv preprint arXiv:2006.01096*, 2020.
  36. Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946*, 2019.
  37. Ajay Kumar Tanwani, Raghav Anand, Joseph E Gonzalez, and Ken Goldberg. Rilaas: Robot inference and learning as a service. *IEEE Robotics and Automation Letters*, 2020.
  38. Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In *2015 IEEE In-*

- formation Theory Workshop (ITW), pages 1–5. IEEE, 2015.
- 39. Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.
  - 40. Chenxi Tu, Eiji Takeuchi, Alexander Carballo, and Kazuya Takeda. Point cloud compression for 3d lidar sensor using recurrent neural network with residual blocks. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 3274–3280. IEEE, 2019.
  - 41. Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(86):2579–2605, 2008.
  - 42. Joshua Vander Hook, Julie Castillo-Rogez, Richard Doyle, Tiago Stegun Vaquero, Trent M Hare, Randolph L Kirk, Valerie Fox, Dmitriy Bekker, and Alice Cocoros. Nebulae: A proposed concept of operation for deep space computing clouds. In *2020 IEEE Aerospace Conference*, pages 1–14. IEEE, 2020.
  - 43. Maurice Weber, Cedric Renggli, Helmut Grabner, and Ce Zhang. Lossy image compression with recurrent neural networks: from human perceived visual quality to classification accuracy. *arXiv preprint arXiv:1910.03472*, 2019.