

# Who2com: Collaborative Perception via Learnable Handshake Communication

Yen-Cheng Liu, Junjiao Tian, Chih-Yao Ma, Nathan Glaser, Chia-Wen Kuo and Zsolt Kira  
Georgia Institute of Technology

{ycliu, jtian73, cyma, nglaser3, albert.cwkuo, zkira}@gatech.edu

*Abstract*—In this paper, we propose the problem of *collaborative perception*, where robots can combine their local observations with those of neighboring agents in a learnable way to improve accuracy on a perception task. Unlike existing work in robotics and multi-agent reinforcement learning, we formulate the problem as one where learned information must be shared across a set of agents *in a bandwidth-sensitive manner* to optimize for scene understanding tasks such as semantic segmentation. Inspired by networking communication protocols, we propose a multi-stage handshake communication mechanism where the neural network can learn to compress relevant information needed for each stage. Specifically, a target agent with degraded sensor data sends a compressed request, the other agents respond with matching scores, and the target agent determines *who to connect with* (i.e., receive information from). We additionally develop the AirSim-CP dataset and metrics based on the AirSim simulator where a group of aerial robots perceive diverse landscapes, such as roads, grasslands, buildings, etc. We show that for the semantic segmentation task, our handshake communication method significantly improves accuracy by approximately 20% over decentralized baselines, and is comparable to centralized ones using a quarter of the bandwidth.

## I. INTRODUCTION

A great deal of progress has been made in single-agent scene understanding using deep neural networks [40, 12, 4]. However, as these methods become ubiquitous and larger numbers of robots are deployed in the world, it becomes beneficial for them to share knowledge via communication. For example, knowledge sharing across a fleet of self-driving cars could alleviate a number of challenges such as occlusion and sensor degradations or failures.

In this paper, we propose the problem of *collaborative perception*, where robots can combine their local observations with those of neighboring agents to improve accuracy in a perception task, such as semantic segmentation [23, 3]. This can result in significant improvements, for example when the receiving agent’s sensors are occluded or degraded (see Figure 1). We therefore formulate a problem where a degraded agent can communicate with other agents to improve its perceptual abilities. Unlike past methods that focus on multi-robot localization and mapping [6, 32], we develop agents that *learn what to communicate* in a manner that is amenable to end-to-end deep learning, which dominates scene understanding. In addition, different from multi-agent reinforcement learning [36, 2, 10, 29, 17, 30, 9], we seek to do so under bandwidth constraints. We therefore propose to learn *who to communicate with* in order to reduce bandwidth

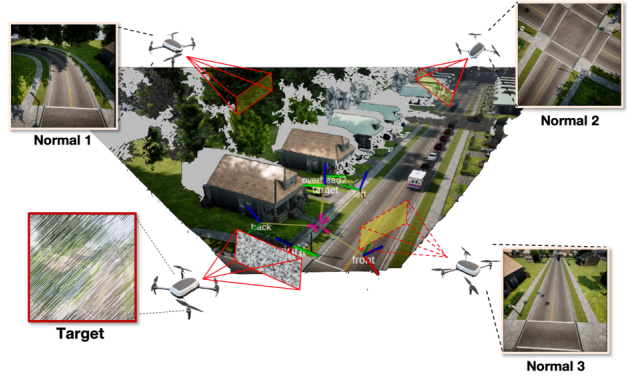


Fig. 1: **Illustration of collaborative perception.** Our collaborative perception task is to improve the perception ability of a degraded agent using information from other agents in a bandwidth-limited way.

requirements while improving accuracy. Such a method is effective for standard observation problems which have the sub-modularity property, i.e., that adding more observation agents achieves diminishing returns [20]. Similar advantages have also been shown using rate distortion theory [8].

In order to investigate the inherent trade-off between accuracy and bandwidth, especially in a manner that scales in a bounded way with respect to the number of agents, we propose a three-stage communication mechanism inspired by three-way handshaking in the regime of communication networking [21]. The three steps of our method are: 1) **request**: the degraded agent broadcasts a compressed request conditioned on its visual observation, 2) **match**: the other agents compute a learned matching score between their own visual observations and the received request, and 3) **connect**: the degraded agent selects one of the agents to communicate with and further improve its own prediction accuracy in downstream perception tasks. The entire mechanism is trained in an end-to-end manner, using only supervision for the down-stream task (e.g., semantic segmentation). We show that this communication mechanism effectively decouples the request, score, and actual transmitted value to allow different (i.e., asymmetric) sizes during this communication. Our experiments demonstrate that this results in significant bandwidth savings when compared to centralized baselines and accuracy gains over uniform compression across agents.

In order to investigate the properties of the resulting

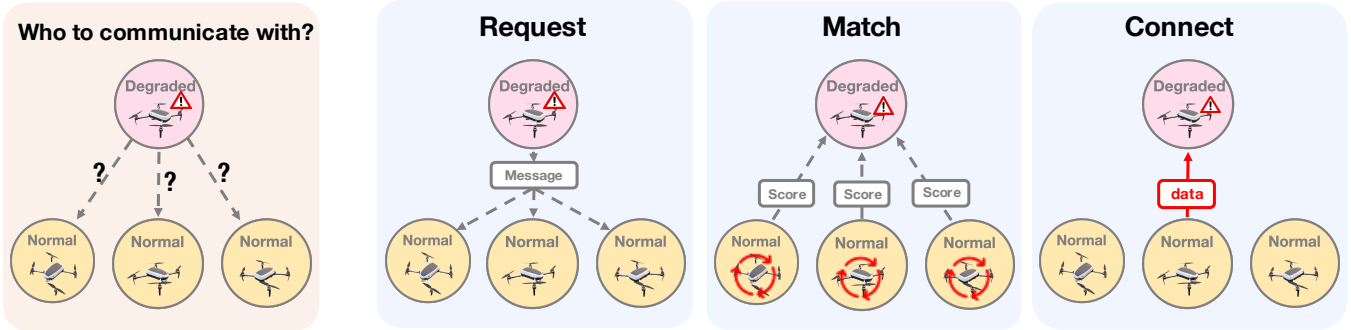


Fig. 2: **The concept of three-stage handshake communication.** Our model consists of three major steps to establish the connection with the selected agent. **Request:** the degraded agent first broadcasts the compressed request message to other agents. **Match:** the normal agents compute matching scores based on the individual observations and the received request message. **Connect:** after returning all scores back to the target degraded agent, it connects to one of the normal agents with the highest matching score and obtains high-bandwidth data (*i.e.*, feature maps for the semantic segmentation task) from it.

method, we develop the AirSim-CP dataset and metrics using the AirSim simulator [34], where a group of aerial robots fly over a map with diverse landscapes, such as roads, grasslands, buildings, lakes, *etc.* We vary a number of elements across the scenarios, including their trajectories, partial or complete overlap between the fields of view of the target and other agents, and with noisy or accurate pose information available for cross-view geometric warping. We quantitatively show that our proposed method is able to significantly improve accuracy by approximately 20% over decentralized baselines, and is comparable to centralized ones using a quarter of the bandwidth.

We highlight the contributions of this paper as follows:

- We propose a new problem in the intersection of multi-agent systems, perception, communication, and learning. Compared with prior works on learning to communicate, we are the first to tackle the problem of learning to communicate with bandwidth constraints to the best of our knowledge.
- Different from other multi-agent systems [14, 17], our collected dataset, AirSim-CP, provides high-resolution and photo-realistic images for better evaluating multi-agent perception tasks with communication.
- We propose an end-to-end communication framework trained without supervision indicating the ground-truth agent for communication, and shows superior accuracy to decentralized baselines and comparably to strong centralized ones with a fraction of the bandwidth.

## II. RELATED WORK

Communication in multi-agent environments is a foundation of both collaborative perception and decision-making. This topic has been studied extensively in the regime of Multi-Agent Reinforcement Learning (MARL) [22, 33, 28]. Early attempts utilized a pre-defined communication protocol to propagate the information across agents [37, 39, 27], while dynamics of environments and a varied number of agents gave rise to the development of learnable multi-agent communication [36, 11, 2, 14, 10, 29, 17, 35, 18, 7, 30, 9, 15, 19]. Existing works on MARL demonstrated the effectiveness of

communication for various tasks, applying the models on simplistic 2D grid environments where the observation of each agent is low-dimensional. As noted in Jain *et al.* [15], studying collaboration in simplistic environments does not permit study of the interplay between perception and communication. Thus, in this work, we examine our framework in a more complicated and photorealistic environment.

Among existing works on learning to communicate, the most relevant framework to our work is TarMac [7], where the targeted communication is defined as the transmitted message determined by both the sender and receiver agents. However, during the communication, both message and data are broadcast to all of the other agents, hence not taking bandwidth usage into account. On the other hand, other recent work proposed to construct the communication group based on either pre-defined rules [17, 16] or a unified communication network [35, 36, 14, 29, 35, 7]. In this way, the bandwidth usage during communication increases as the number of agents scales up. In contrast, our framework aims to minimize the bandwidth consumption yet maintain performance for the perception task by selecting which agent(s) to communicate with.

## III. PROPOSED METHOD

### A. Problem Definition and Motivation

In our proposed *collaborative perception* task, an environment consists of  $N$  agents with their own observations  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1,\dots,N}$ , while the observation of the target agent  $\tilde{\mathbf{x}}_j, j \in \{1, \dots, N\}$  is degraded. The goal of the target agent is to integrate information received from other agents to derive the prediction for its own local observation. Realistically, communication mechanisms often have bandwidth limitations, preventing the transmission of a large quantity of information during communication. Thus, our goal is to derive a distributed and information-fusing framework which is able to (1) maximize the prediction accuracy of the downstream perception tasks for the target agent and (2) minimize the bandwidth used during transmission, as illustrated in Figure 1. We propose a communication framework that can

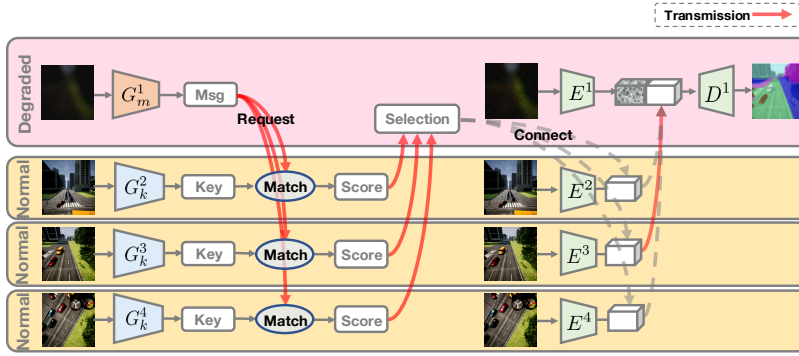


Fig. 3: **Our proposed model and communication steps.** The degraded agent first computes a low-dimensional vector as the request message and broadcasts it to other normal agents, and each normal agent then generates a matching score conditioned on the request message. Finally, the degraded agent accesses the requested information from one of the agents during inference (and from all agents during training). The proposed method is trained end-to-end and significantly improves accuracy over decentralized baselines while minimizing bandwidth usage over centralized ones.

be generalized to many perception tasks, and in this paper we consider semantic segmentation as one instantiation to evaluate the framework.

### B. Communication via Three-Stage Handshake

Rather than broadcasting all information within the entire network in a brute-force manner, an efficient way of minimizing the bandwidth usage while maintaining performance is to select which agent to communicate with. Toward this end, inspired by related work in communication networks [21], we introduce a three-stage handshake communication mechanism shown in Figure 2. Such a selection is motivated by works on rate distortion theory [8] and sub-modularity of multiple observations [20]. Our empirical experimental results also show that this approach yields a better trade-off between bandwidth usage and performance, compared to uniform compression of all information from the agents.

Our communication mechanism consists of three major steps: **request**, **match**, and **connect**. Specifically, the degraded agent first broadcasts its request message,  $\mu_j \in \mathbb{R}^m$ , to the neighboring normal agents, and the normal agents compute a matching score,  $s_{ji}$ , between their keys,  $\kappa_i \in \mathbb{R}^k$ , and the request message. Once the normal agents return their matching scores back to the degraded agent, the degraded agent further selects the best  $n$  agents to connect with (*i.e.*, receive information from) according to these matching scores<sup>1</sup>. The complete communication steps are illustrated in Figure 3. In each step, information can be compressed by each agent  $i$  through a key generator  $G_k^i$ , a message generator  $G_m^i$ , an image encoder  $E^i$ , and an task decoder  $D^i$ . Note that this approach effectively decouples the various stages, allowing for different compression rates for the message, keys, and values. This significantly benefits the trade-off between bandwidth and accuracy. We now detail the steps:

**Request** - The degraded agent  $j$  first compresses its

observation  $\tilde{x}_j$  to a low-dimensional message  $\mu_j$ :

$$\mu_j = G_m^j(\tilde{x}_j; \theta_m) \in \mathbb{R}^m, \quad (1)$$

where  $G_m^j$  is a message generator parameterized by  $\theta_m$ . The propagated message  $\mu_j$  compresses important information from the local observation of the degraded agent  $j$ .

**Match** - In the match step, each of other agents derives the matching score  $s_{ji}$  between the received request message  $\mu_j$  from the degraded agent and the key  $\kappa_i$  generated from its own observations,

$$s_{ji} = \Phi(\mu_j, \kappa_i), \quad \kappa_i = G_k^i(x_i; \theta_k) \in \mathbb{R}^k, \quad (2)$$

where  $\Phi(\cdot, \cdot)$  represents the matching function of two vectors and  $G_k^i$  denotes a key generator parameterized by  $\theta_k$ . We use the *general* attention mechanism [25] as our matching function:

$$\text{General: } \Phi = \mu_j^T \mathbf{W}_a \kappa_i, \quad (3)$$

where  $\mathbf{W}_a$  is a learnable parameter. We also compare it with two other attention mechanisms: Scale Dot-Product  $\Phi = \mu_j^T \kappa_i / \sqrt{d_n}$  [38] and Additive  $\Phi = \mathbf{W}_a^T \tanh(\mathbf{W}_k \kappa_i + \mathbf{W}_m \mu_j)$  [1], where  $\mathbf{W}_k, \mathbf{W}_m$  denote parameters to be learned and  $d_n$  denotes the dimension of the message and keys. Note that only the general and additive functions allow for different key and message sizes. The scale dot-product function requires identical message-key size. Empirically, we find that *general* attention achieve the best performance in our experiments and hence use it unless otherwise specified. Note that we assume equal cost links between agents, though our models can further support per-link costs (unlike the baseline methods). We leave this for future work.

**Connect** - In the connect step, the selected  $\hat{i}$ -th agent transmits the requested information (*e.g.*, a feature map for semantic segmentation)  $f_{\hat{i}}$  to the degraded agent. With the integrated information, the target degraded agent makes a final prediction  $\tilde{y}_j$  as follows:

$$\tilde{y}_j = D^j([\tilde{f}_j; f_{\hat{i}}; \theta_d]), \quad (4)$$

<sup>1</sup>Note that in this paper we select the top  $n = 1$ , but the method can be generalized to top- $n$  selection.

where  $\mathbf{f}_i = E^i(\mathbf{x}_i; \theta_e) \in \mathbb{R}^{d_f \times d_f \times d_c}$  is the feature map from the local observation of normal agent  $i$ ,  $d_f, d_c$  are the spatial dimension and number of channels of the feature maps respectively,  $\tilde{\mathbf{f}}_j = E^j(\tilde{\mathbf{x}}_j; \theta_e)$  is the feature map from the noisy observation of the degraded agent, and  $[\cdot; \cdot]$  is the concatenation operator along the channel dimension.  $\theta_e$  and  $\theta_d$  are the parameters in the task encoder and decoder.

### C. Learning to Communicate with Weak Supervision

**Centralized training with decentralized execution.** Our learning strategy is inspired by the concept of centralized training with decentralized execution [24]. During training our target agent can access the observations of all other agents. On the other hand, during inference the target agent is required to perform in a bandwidth-limited manner by only accessing information from the selected agent(s).

Specifically, during training, the task decoder uses the sum of observations from all normal agents weighted by their corresponding matching scores and further computes the final prediction akin to Eq. 4 during inference:

$$\tilde{\mathbf{y}}_j = D^j([\tilde{\mathbf{f}}_j; \mathbf{f}_{sum}]; \theta_d), \quad \mathbf{f}_{sum} = \sum_{i=1}^N \alpha_{j,i} \mathbf{f}_i, \quad (5)$$

where  $\alpha_{j,i}$  is  $i$ th element of  $\alpha_j = \rho([s_{j1}; \dots; s_{jN}]) \in \mathbb{R}^N$  and  $\rho$  is a softmax operation. The most straightforward decentralized execution method is simply adopting *argmax* selection, *i.e.*, connecting only to the agent with the highest computed matching score:

$$\hat{i} = \underset{i}{\operatorname{argmax}} s_{ji}. \quad (6)$$

However, *argmax* selection is non-differential during training. We address this issue by simply applying the *softmax* operator in the training stage and *argmax* operator in the inference stage. Empirically, we find that this simple method achieves similar results compared to other more complex schemes, *e.g.*, *sparsemax* [26]. Note also that this can be generalized to top- $n$  selection as well.

**Training objective.** Learning our model does not require supervision of ground truth labels indicating the best agent to communicate with. Therefore, the only supervision for our model comes from the ground truth annotation at the target view (*e.g.*, segmentation masks). The objective function of our model, which is trained end-to-end, can thus be defined as  $\mathcal{L} = \mathcal{H}(\mathbf{y}_j, \tilde{\mathbf{y}}_j)$ , where  $\mathcal{H}$  is the cross-entropy loss, and  $\mathbf{y}_j$  denotes the ground truth labels of the target agent’s view.

## IV. EXPERIMENTAL RESULTS

### A. AirSim-CP Dataset

**Dataset.** Our AirSim-CP dataset is built upon the AirSim simulator [34], where a group of **five** drones fly over a map with diverse landscapes, such as roads, grasslands, buildings, lakes, *etc.* Currently, in our AirSim-CP dataset, we use semantic segmentation as the downstream task to benchmark methods for the collaborative perception problem. For each drone, RGB images, depth images, and poses are recorded.

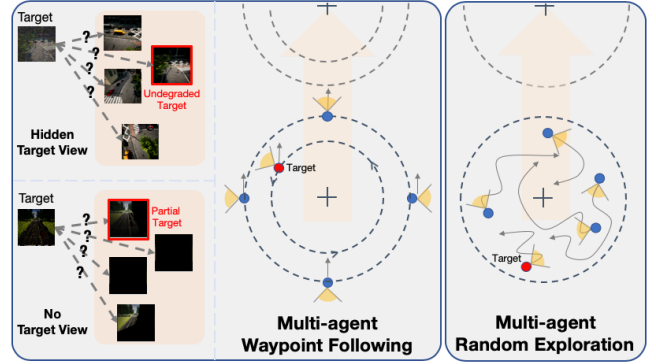


Fig. 4: Illustration of our experimental cases.

We also provide the semantic segmentation mask of one of agents.

### B. Proposed Experimental Settings.

As illustrated in Figure 4, in order to obtain perceptual data under realistic trajectories, the drones perform pre-specified tasks of waypoint following and multi-agent random exploration. We then consider **four** experimental settings: (1) hidden target view (waypoint following trajectories), (2-3) accurate or inaccurate pose (waypoint following trajectories), and (4) accurate pose (random exploration trajectories). We collect approximately 10-20k images per setting, with a roughly 60%/20%/20% train/val/test split. In all cases, there is a degraded target agent. We perturb the target agent’s view by applying a Gaussian blur filter with random size (from 1 to 100) and Gaussian noise. When specified, depth and pose information is used to warp the normal agents’ views to the target view. In terms of the optimization, we use ResNet18 [13] as our feature backbone and train it for 200k iterations with Adam optimizer with learning rate  $10^{-5}$ .

**Hidden Target View with Multi-agent Waypoint Following.** The agents are asked to navigate between waypoints but when performing the semantic segmentation task, rather than the *normal* agents being all of the neighboring agents, we replace one of the normal agent with an *un-degraded* version of the target view. This task tests whether the baseline and the proposed method can help the target agent find the *hidden* ground-truth target view in one of the neighboring agents. This experiment thus can be regarded as a sanity check on communication. Note that we do not use any 3D information to warp normal views to the target view. The motivation for this case is to remove the confound of geometric warping and image misalignments (which is important for the semantic segmentation task [31]) from the study of the collaborative perception task. Our focus is to make sure that communication can be effective and accurate in a bandwidth-limited manner. Hence, we can directly assess only communication effectiveness from this case.

**Accurate Pose with Multi-agent Waypoint Following.** Similar to the previous setting, the five drones are performing waypoint following jointly. Differently, the field-of-view (FOV) of the target agent only partially overlaps

with some subset of the normal agents. This case is designed to test whether the proposed method can select the normal agent with partially overlapping FOVs, in order to aid the downstream perception task. In order to maintain the image alignment for a better segmentation prediction, we use 3D information from the depth maps of each normal agent’s view and an accurate relative pose transformation to the target view to warp the pixels of the normal agent’s observations to the target’s view. Note that the depth maps do not have to be transmitted (warping is done locally on each agent) but we include the transmission of the target’s pose to the other agents in the bandwidth used, although it is minimal.

**Inaccurate Pose with Multi-agent Waypoint Following.**

To further make our experimental setting more realistic, we add noise in the agents’ positions. This results in warped images that are not well-aligned with the target view.

**Accurate Pose with Multi-agent Random Exploration.**

We also investigate collaborative perception during multi-agent random exploration, where agents approach a target location, disperse, and wander. As the agents explore the environment individually, the relative positions and overlapping fields-of-view of agents change frequently.

*C. Baselines*

We consider the following methods for comparison:

- *Single normal (upper bound) and Single degraded (lower bound)* : the models are trained with single non-degraded and degraded images respectively for the target agent.
- *CatAll (centralized)*: the model uses the concatenation of all features from both degraded and normal agents as input for semantic segmentation.
- *Attention (centralized)*: the *Attention* mechanism weights and sums up feature maps instead of the concatenation of the *CatAll* method.
- *Compression (centralized)*: the compression model applies two additional convolutional layers and performs uniform compression of all observations at rate 25%, with concatenation used for combining them. Note that we can certainly replace our image encoder with more sophisticated compression encoders to further improve the compression rate [5].
- *Random selection (distributed)*: instead of learning to select which agent to communicate with, here the feature map from a random normal agent is selected.
- *Ours (distributed)*: we denote our proposed method as ours with message (*ours w/ msg*), and another variant where the message request  $\mu_j$  is set to a constant vector with ones to check whether the message request is essential. It is worth noting that we *do not* use any label indicating the best agent during the training.

Both *CatAll* and *Attention* require all feature maps from normal agents to be sent to the degraded agent. The bandwidth of the centralized baselines scales linearly with the number of agents in the system, while *Random selection* and *ours* only requires a single image feature map to be transmitted.

*D. Evaluation metrics.*

In order to evaluate and analyze the effectiveness of models, we use 1) overall accuracy to measure the performance of semantic segmentation [23] and 2) Kbytes per frame (kbpf) to measure the bandwidth usage (BW) to examine the ability of communication and selection. In addition, to better benchmark the performances on the collaborative perception task under limited bandwidth, we introduce the **Bandwidth-Improvement Score** (BIS) defined as:

$$BIS = \frac{\delta - \bar{\delta}}{(\hat{\delta} - \bar{\delta})\omega}, \tag{7}$$

where  $\delta$  is the overall accuracy of the examined method,  $\bar{\delta}$  is the overall accuracy of the single degraded model (*i.e.*, lower bound on overall accuracy),  $\hat{\delta}$  is the overall accuracy of the single normal model (*i.e.*, upper bound on overall accuracy), and  $\omega$  is the bandwidth usage (in Mbytes per frame) of the examined method. The BIS score is defined as a ratio of relative improvement in overall accuracy over bandwidth usage. Smaller bandwidth usage and larger improvements in overall accuracy lead to higher scores.

V. RESULTS AND ABLATION STUDIES

In this section, we compare our proposed method with baselines on the four cases (Sec. IV-B) as shown in Table I.

**Hidden Target View with Multi-agent Waypoint Following.** As mentioned, geometric warping is not applied for this case. This results in better evaluation of communication, since the warping noise for the target view is removed and non-warped normal images make the selection more difficult.

Several observations can be made from this case. First, as the centralized baselines are able to access all observations from different views, it should be the upper bounds of all distributed methods including our proposed models. However, we find that our model improves overall accuracy by a relative 16.51% with respect to *CatAll* with only one quarter bandwidth usage. This shows that simply concatenating all of the information cannot guarantee that the network will combine it meaningfully and bandwidth is likely to be wasted. Second, in order to predict pixel-wise outputs and accurately predict fine-grained classes, scene understanding tasks require high-dimensional feature maps during inference. Using the overly compressed feature map may degrade the overall accuracy, hence our model with the message can improve overall accuracy a relative 20.06% with respect to the *Compression* model. Lastly, we also observe that our model with the message can improve mIoU by a relative 29.49% compared to without the message. This demonstrates the necessity of the message in the communication.

**Accurate and Inaccurate Pose with Multi-agent Waypoint Following.** One challenge for these cases is that the field of view between target and normal agents are partially overlapping. With the forward geometric warping on the visual observations of the normal agents, we observe that only one or two agents contain(s) partial information of the target view. Thus the performances of distributed models

TABLE I: Experimental results on Multi-agent Waypoint Following and Random Exploration.

		Waypoint Following (Hidden Target View)		Waypoint Following (Accurate Pose)		Waypoint Following (Inaccurate Pose)		Random Exploration (Accurate Pose)		
		BW (kpbf)	Overall Acc	BIS	Overall Acc	BIS	Overall Acc	BIS	Overall Acc	BIS
Upper bound	Single Normal	-	88.14	-	88.13	-	89.7	-	89.16	-
Centralized	CatAll	4096	72.58	0.049	80.05	0.17	80.33	0.172	78.74	0.133
	Attention	4096.03	69.08	0.004	84.38	0.212	82.79	0.174	81.07	0.159
	Compression	1024	70.44	0.085	76.2	0.527	76.93	0.441	73.05	0.277
Distributed	Random Selection	1024	69.16	0.019	64.94	0.082	67.49	0.028	69.58	0.123
	Ours w/o msg	1024.03	65.31	-0.179	78.10	0.602	79.34	0.546	80.69	0.621
	Ours w/ msg	1028.03	84.57	<b>0.812</b>	80.42	<b>0.691</b>	79.44	<b>0.549</b>	80.97	<b>0.631</b>
Lower bound	Single Degraded	-	68.79	-	62.88	-	66.84	-	66.85	-

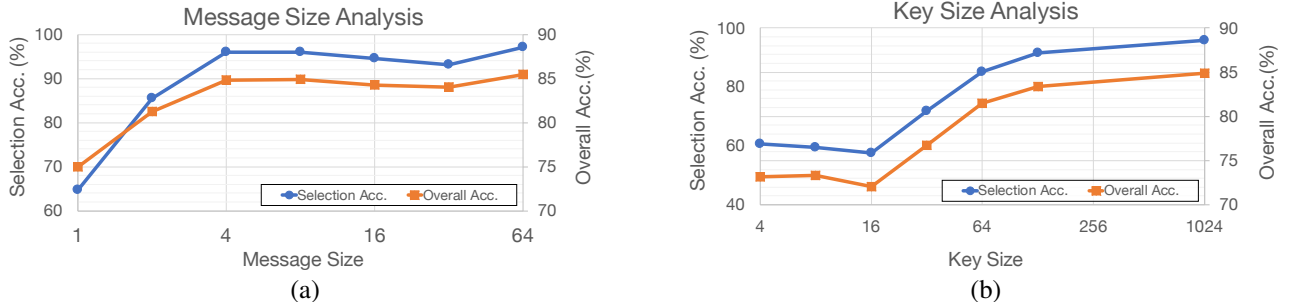


Fig. 5: Ablation study on varying (a) message and (b) key size. We vary the size of key from 4 to 1024 and vary the size of message from 1 to 64 on Hidden Target View. Note that we use the key with size 8 on the varying message analysis, while the message with size 1024 is used for key size analysis. Without the message, the model performs the selection accuracy of 25.52% and the overall accuracy of 61.4%, and we use key size of 1024 for (a) and message size of 8 for (b).

drop compared to the previous task. Our models are still able to achieve similar results compared with centralized methods, using only a quarter of the bandwidth, and still achieve the highest BIS among these methods. It is worth noting that our models with and without the message perform similarly because the model can rely on cues from warping (*e.g.*, amount of overlap) from the normal agent images, hence reducing the need for conditioning on the message. Note also that although the BIS degrades somewhat due to inaccurate pose, it is still able to significantly beat the baselines.

#### Accurate Pose with Multi-agent Random Exploration.

Although overlapping FOVs and motions of agents change frequently in this case, we observe that our models still perform favorably against other baseline methods. This shows our models’ robustness in different environments and tasks.

**Message and key sizes.** To better examine the accuracy of agent selection, we manually label the “best” agent of the testing sets of Hidden Target View and further measure the selection accuracy and overall accuracy of various models.

When using the *general* attention, the size of message and key can be set to different values. We first analyze the effect of message size by varying it from 1 to 64 as shown in Figure 5a. We observe that a larger message size results in increased selection accuracy and segmentation quality. Note that a message size of 4 is sufficient to achieve comparable performance on both selection and segmentation, after which the performance plateaus. We also conduct a similar experiment by varying the key size from 4 to 1024, and the same trend can be observed in Figure 5b. Also, we empirically

found that a small message (*e.g.*, 8) paired with large key size (*e.g.*, 1024) can achieve amenable performance. This is also the sizes of our models for all experiments in Sec. V. Importantly, the effectiveness of asymmetric sizes is an interesting finding as it allows us to use a small size for the message that is sent to other agents while using larger sizes for the key used locally by each agent to compute the scores (and hence does not need to be transmitted). These results validate this advantage.

## VI. CONCLUSION

In this paper, we formulated the problem of *collaborative perception*, where agents can combine their local observations with those of other agents in order to improve performance on scene understanding tasks. Inspired by the network communication literature, we propose a handshake communication mechanism for which the network can learn compressed representations. Key to our approach is that we decouple the message, key, and value elements to support asymmetric compression, resulting in bandwidth savings. We introduce the AirSim-CP dataset and benchmarking metrics to evaluate our method, and show that our method is able to effectively combine information from neighboring agents to improve accuracy using significantly less bandwidth than centralized approaches.

## VII. ACKNOWLEDGEMENT

This work was supported by ONR grant N00014-18-1-2829.

## REFERENCES

- [1] Dzmitry Bahdanau, KyungHyun Cho, and Yoshua Bengio. “Neural machine translation by jointly learning to align and translation”. In: *Proceedings of the International Conference on Learning Representations (ICLR)*. 2015.
- [2] Peter Battaglia et al. “Interaction networks for learning about objects, relations and physics”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2016.
- [3] Liang-Chieh Chen et al. “Rethinking atrous convolution for semantic image segmentation”. In: *arXiv preprint arXiv:1706.05587* (2017).
- [4] Po-Yi Chen et al. “Towards Scene Understanding: Unsupervised Monocular Depth Estimation With Semantic-Aware Representation”. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2019.
- [5] Zhengxue Cheng et al. “Learning Image and Video Compression through Spatial-Temporal Energy Compaction”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019, pp. 10071–10080.
- [6] Alexander Cunningham, Manohar Paluri, and Frank Dellaert. “DDF-SAM: Fully distributed SLAM using constrained factor graphs”. In: *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE. 2010, pp. 3025–3030.
- [7] Abhishek Das et al. “TarMAC: Targeted Multi-Agent Communication”. In: *Proceedings of the International Conference on Machine Learning (ICML)*. 2019.
- [8] Roel Dobbe, David Fridovich-Keil, and Claire Tomlin. “Fully decentralized policies for multi-agent systems: An information theoretic approach”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2017, pp. 2941–2950.
- [9] Jakob N Foerster et al. “Counterfactual multi-agent policy gradients”. In: *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*. 2018.
- [10] Jakob Foerster et al. “Learning to communicate with deep multi-agent reinforcement learning”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2016.
- [11] Jakob Foerster et al. “Learning to communicate with deep multi-agent reinforcement learning”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2016.
- [12] Clément Godard, Oisín Mac Aodha, and Gabriel J Brostow. “Unsupervised monocular depth estimation with left-right consistency”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2017, pp. 270–279.
- [13] Kaiming He et al. “Deep residual learning for image recognition”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [14] Yedid Hoshen. “Vain: Attentional multi-agent predictive modeling”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2017.
- [15] Unnat Jain et al. “Two Body Problem: Collaborative Visual Task Completion”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019.
- [16] Jiechuan Jiang, Chen Dun, and Zongqing Lu. “Graph Convolutional Reinforcement Learning for Multi-Agent Cooperation”. In: *arXiv preprint arXiv:1810.09202* (2018).
- [17] Jiechuan Jiang and Zongqing Lu. “Learning attentional communication for multi-agent cooperation”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2018.
- [18] Daewoo Kim et al. “Learning to Schedule Communication in Multi-agent Reinforcement Learning”. In: *Proceedings of the International Conference on Learning Representations (ICLR)*. 2019.
- [19] Woojun Kim, Myungsik Cho, and Youngchul Sung. “Message-Dropout: An Efficient Training Method for Multi-Agent Deep Reinforcement Learning”. In: *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*. 2019.
- [20] Andreas Krause and Carlos Guestrin. “Near-optimal observation selection using submodular functions”. In: *AAAI*. Vol. 7. 2007, pp. 1650–1654.
- [21] James F Kurose. *Computer networking: A top-down approach featuring the internet, 3/E*. Pearson Education India, 2005.
- [22] Michael L Littman. “Markov games as a framework for multi-agent reinforcement learning”. In: *Machine learning proceedings 1994*. Elsevier, 1994, pp. 157–163.
- [23] Jonathan Long, Evan Shelhamer, and Trevor Darrell. “Fully convolutional networks for semantic segmentation”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2015, pp. 3431–3440.
- [24] Ryan Lowe et al. “Multi-agent actor-critic for mixed cooperative-competitive environments”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2017.
- [25] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. “Effective approaches to attention-based neural machine translation”. In: *Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2015.
- [26] Andre Martins and Ramon Astudillo. “From softmax to sparsemax: A sparse model of attention and multi-label classification”. In: *Proceedings of the International Conference on Machine Learning (ICML)*. 2016, pp. 1614–1623.
- [27] Francisco S Melo, Matthijs TJ Spaan, and Stefan J Witwicki. “QueryPOMDP: POMDP-based communication in multiagent systems”. In: *Proceedings of*

*the 9th European conference on Multi-Agent Systems*. 2011.

- [28] Liviu Panait and Sean Luke. “Cooperative multi-agent learning: The state of the art”. In: *Autonomous agents and multi-agent systems* (2005).
- [29] Peng Peng et al. “Multiagent bidirectionally-coordinated nets: Emergence of human-level coordination in learning to play starcraft combat games”. In: *arXiv preprint arXiv:1703.10069* (2017).
- [30] Emanuele Pesce and Giovanni Montana. “Improving Coordination in Multi-Agent Deep Reinforcement Learning through Memory-driven Communication”. In: *Neural Information Processing Systems Workshops*. 2019.
- [31] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. “U-net: Convolutional networks for biomedical image segmentation”. In: *International Conference on Medical image computing and computer-assisted intervention*. Springer. 2015, pp. 234–241.
- [32] Sajad Saeedi et al. “Multiple-robot simultaneous localization and mapping: A review”. In: *Journal of Field Robotics* 33.1 (2016), pp. 3–46.
- [33] Jurgen Schmidhuber. “A general method for multi-agent reinforcement learning in unrestricted environments”. In: *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*. 1996.
- [34] Shital Shah et al. “AirSim: High-Fidelity Visual and Physical Simulation for Autonomous Vehicles”. In: *Field and Service Robotics*. 2017.
- [35] Amanpreet Singh, Tushar Jain, and Sainbayar Sukhbaatar. “Learning when to Communicate at Scale in Multiagent Cooperative and Competitive Tasks”. In: *arXiv preprint arXiv:1812.09755* (2018).
- [36] Sainbayar Sukhbaatar, Rob Fergus, et al. “Learning multiagent communication with backpropagation”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2016.
- [37] Ming Tan. “Multi-Agent Reinforcement Learning: Independent vs. Cooperative Agents”. In: *Proceedings of the International Conference on Machine Learning (ICML)*. 1993.
- [38] Ashish Vaswani et al. “Attention is all you need”. In: *Advances in Neural Information Processing Systems (NIPS)*. 2017.
- [39] Chongjie Zhang and Victor Lesser. “Coordinating multi-agent reinforcement learning with limited communication”. In: *Proceedings of International Conference on Autonomous agents and Multi-agent Systems (AAMAS)*. 2013.
- [40] Hengshuang Zhao et al. “Pyramid Scene Parsing Network”. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. July 2017.