

UMC: A Unified Bandwidth-efficient and Multi-resolution based Collaborative Perception Framework

Tianhang Wang¹, Guang Chen^{1,*}, Kai Chen¹, Zhengfa Liu¹, Bo Zhang², Alois Knoll³, Changjun Jiang¹
¹Tongji University, ²Westwell lab, ³Technische Universität München
 {tianya.wang, guangchen, 14sdck, cjjiang}@tongji.edu.cn
 zhengfaliu2011@163.com, zhangaigh@gmail.com, knoll@in.tum.de

Abstract

Multi-agent collaborative perception (MCP) has recently attracted much attention. It includes three key processes: communication for sharing, collaboration for integration, and reconstruction for different downstream tasks. Existing methods pursue designing the collaboration process alone, ignoring their intrinsic interactions and resulting in suboptimal performance. In contrast, we aim to propose a Unified Collaborative perception framework named UMC, optimizing the communication, collaboration, and reconstruction processes with the Multi-resolution technique. The communication introduces a novel trainable multi-resolution and selective-region (MRSR) mechanism, achieving higher quality and lower bandwidth. Then, a graph-based collaboration is proposed, conducting on each resolution to adapt the MRSR. Finally, the reconstruction integrates the multi-resolution collaborative features for downstream tasks. Since the general metric can not reflect the performance enhancement brought by MCP systematically, we introduce a brand-new evaluation metric that evaluates the MCP from different perspectives. To verify our algorithm, we conducted experiments on the V2X-Sim and OPV2V datasets. Our quantitative and qualitative experiments prove that the proposed UMC greatly outperforms the state-of-the-art collaborative perception approaches.

1. Introduction

Single-vehicle perception has made remarkable achievements in object detection[21, 31, 32], segmentation[29, 33], and other tasks with the advent of deep learning. However, single-vehicle perception often suffers from environmental conditions such as occlusion[39, 50] and severe weather[4, 16, 53], making accurate recognition challenging. To overcome these issues, several appealing studies have been devoted to collaborative perception[1, 2, 7, 42, 49], which take

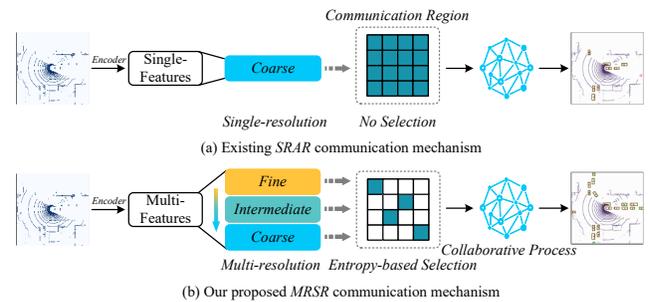


Figure 1. **The SRAR versus the proposed MRSR mechanism.** (a) The existing single-resolution and all-region (SRAR) method utilizes the coarse-grained feature with bandwidth inefficient all-region transmission. (b) The proposed MRSR mechanism incorporates bandwidth-efficient entropy-based selection with multi-resolution features.

advantage of sharing the multiple-viewpoint of the same scene with the Vehicle-to-Vehicle(V2V) communication[3].

To design a collaborative perception algorithm, current approaches[40, 24, 45, 19, 44] mainly focus on the collaboration process alone, aiming to design a high-performance collaboration strategy. Nevertheless, such a strategy ignores intrinsic interactions with two critical processes: communication[23, 12] and reconstruction[9]. This kind of collaboration strategy will inherently cause suboptimal performance because the quality of communication directly determines the performance of collaboration[19], and the reconstruction influences the quality of feature maps for downstream tasks. Meanwhile, the collaboration can also affect the design of communication and reconstruction. Hence, failing to optimize these three processes will degrade the system.

In this paper, to the best of our knowledge, we are the first to propose a unified collaborative perception framework that optimizes the communication, collaboration, and reconstruction processes with multi-resolution technique. As for communication, as shown in Figure 1, we propose a

*Corresponding author: guangchen@tongji.edu.cn

novel multi-resolution and selective-region (MRSR) mechanism, which is different from the single-resolution and all-region (SRAR) mechanism that is used by existing collaborative algorithms[19, 45, 40]. In MRSR, we replace the coarse-grained feature in SRAR with multi-grain features, which reflect the scene from a global to a local perspective. The multi-grain features will have the complementary advantages of global structure information and local texture details. Meanwhile, to reduce the burden of bandwidth, unlike SRAR, which blindly transmits all regions, we implement a trainable and region-wise entropy-based communication selection (Entropy-CS) module that highlights the informative regions and selects appropriate regions for each level of resolution. The proposed MRSR is adaptive to real-time communication, reflecting dynamic connections among agents.

We propose a graph-based collaborative GRU(G-CGRU) module for MRSR’s each level resolution feature in collaboration. Compared to the existing collaboration strategies that only focus on the present shared information without considering the previous state, we redesign the GRU-cell[5] to adapt the time series of collaborative perception, which models the continuity of vehicle movement. In G-CGRU, the collaboration depends not only on collaborators but also on the last moment information of the ego agent. Meanwhile, we propose matrix-valued gates in G-CGRU to ensure collaboration at a high spatial resolution and allow the agents to adaptively weight the informative regions.

To adapt the proposed multi-resolution collaboration mechanism for reconstruction, we propose a multi-grain feature enhancement (MGFE) module to strengthen the feature reconstruction process. In MGFE, the fine-grained collaborative feature maps will give direct guidance to the agents’ feature maps via a global pooling operation, and the coarse-grained collaborative feature maps will again enhance the agents’ feature maps from a global perspective. This design allows the agents to comprehensively reconstruct the feature maps from local and global viewpoints.

The general evaluation metric (*e.g.*, *average precision*) can reflect the overall performance of multi-agent collaborative perception (MCP). However, the performance enhancement brought by the collaboration concept that MCP introduced cannot be reflected directly. To address this issue, we introduce a brand new evaluation metric that systematically evaluates the MCP from four aspects.

To validate the proposed framework, we conducted comprehensive experiments in 3D object detection on V2X-Sim[17] and OPV2V[46] datasets. Our proposed unified, bandwidth-efficient and multi-resolution based collaborative perception framework (UMC) achieves a better performance-bandwidth trade-off than the state-of-the-art SRAR-based collaboration methods. Our contributions are listed as follows:

- We present a unified, bandwidth-efficient framework (UMC) for collaborative perception, which optimizes the communication, collaboration, and reconstruction processes with the multi-resolution technique.
- We propose a novel multi-resolution and selective-region mechanism for communication and the graph-based collaborative GRU for each resolution collaboration and multi-grain feature enhancement module for reconstruction.
- We present a brand new evaluation metric for 3D object collaborative perception, which can evaluate the performance from different perspectives.

2. Related works

2.1. V2X Collaborative Perception

F-Cooper[1] is the first to apply feature-level collaboration, which weights interaction information equally using a max-based function. When2com[23] employs an attention mechanism to build a bandwidth-efficient communication group. V2VNet[40] proposes a spatially-aware message transmission mechanism that assigns varying weights to the different agents. DiscoNet[19] adapts the knowledge distillation to enhance student model performance by constraining the teacher model, which is based on corresponding raw data collaboration. V2X-ViT[45] presents a unified V2X framework based on Transformer that takes into account the heterogeneity of V2X systems. Where2comm[12] takes the advantage of sparsity of foreground information in detection downstream task to save the communication bandwidth of V2X collaboration. CoBEVT[44] generates BEV map predictions with multi-camera perception by V2X.

2.2. 3D Perception for Autonomous Driving

Due to the declining price of commodity LiDAR sensors, 3D perception has become more prevalent in autonomous driving. VoxelNet[54] pioneered utilizing 3D object detection in autonomous driving. More and more advanced algorithms have been discovered. There are three branches of 3D perception input forms that have been widely adopted: point-based methods[29, 37, 30] directly process raw captured point; and voxel-based methods[54, 47, 13, 36] divide the points into a volumetric grid; and 2D projections based methods[38, 48, 14] compact the input by projecting points into images. In UMC, we utilize bird’s eye view projection as input, which decreases the processing time and enables the use of 2D object detection networks.

2.3. Multi-resolution Image Inpainting

Multi-resolution is a widely used image inpainting[51, 28, 41] technique to combine low-resolution feature

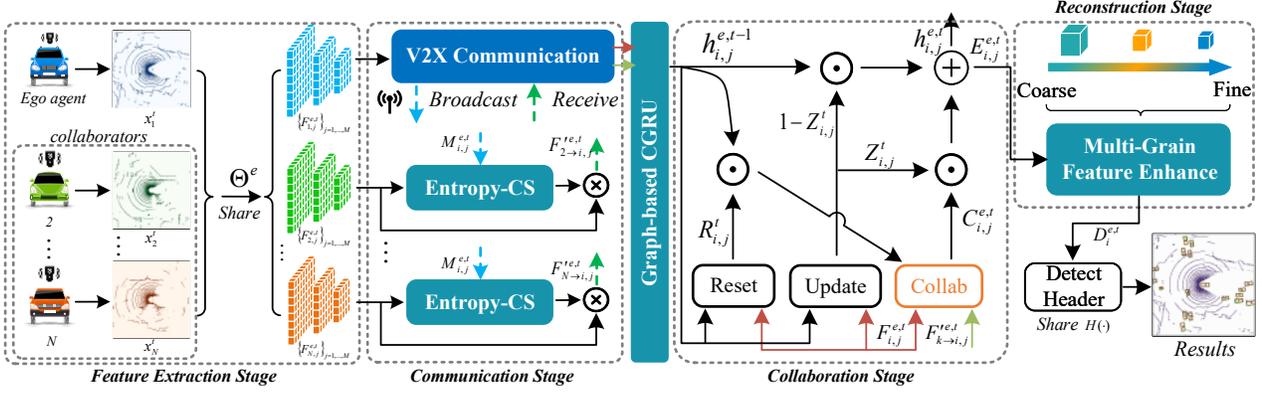


Figure 2. The overview of the proposed UMC framework. i) **Feature extraction stage**, the agents obtain the $F_i^{e,t}$ by the shared feature encoder Θ^e with the observation x_i^t . ii) **Communication stage**, the ego agent (in blue) will broadcast compact query matrix $M_i^{e,t}$ in each resolution to collaborators by V2X communication, and the collaborators will compute the transmission map by entropy-CS module at local. Then the ego agents will receive the selected messages from the assisted collaborators. iii) **Collaboration stage**, the ego agents will employ the G-CGRU module in each resolution for high efficient collaboration. iv) **Reconstruction stage**, the MGFE module will reconstruct the ego agent’s feature by multi-resolution collaborative feature maps for different downstream tasks.

maps that contain global structure information with high-resolution feature maps that are rich in local texture details, such as, PEN-Net[52] proposes the attention map mechanism, which learns from high-resolution feature maps to guide low-resolution features. Wang *et al.* [41] design a parallel multi-resolution fusion network that can effectively fuse structure and texture information. DFNet[11] employs U-Net[34] based multiple fusion blocks to implement multi-resolution constraints for inpainting. In UMC, we implement the multi-resolution technique for a new application scenario: multi-agent collaborative perception. The multi-grain collaborative feature makes it possible for the agents to reconstruct the occlusion area from both a local and a global point of view.

3. UMC

In this section, we present the technical details of the proposed UMC. The overall architecture of the proposed UMC is shown in Figure 2, which consists of three new modules, entropy-CS, G-CGRU, and the MGFE. Entropy-CS extends the traditional entropy theory to select the informative regions of observations, reducing the heavy bandwidth burden brought by multi-resolution. G-CGRU redesigns the GRU-Cell to adapt the multi-agent collaboration, which models the continuity of vehicle movement. MGFE introduces the multi-grain feature enhancement to strengthen the feature reconstruction process for different downstream tasks.

3.1. Problem Statement and Notation

We assume there are N collaborators with their corresponding observations $\mathbf{X}^t = \{x_n^t\}_{n=1,\dots,N}$ at time t . Among those collaborators, the feature set of the i -th agent is defined as $F_i^{e,t} = \{F_{i,j}^{e,t}\}_{j=1,\dots,M} \leftarrow \Theta^e(x_{i,t})$, where

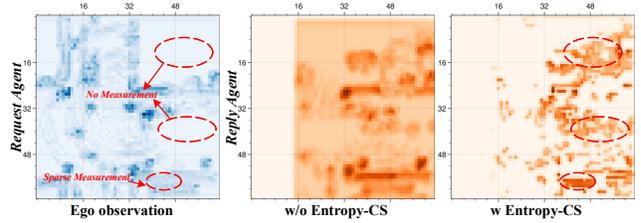


Figure 3. Illustration of entropy-based communication selection. The entropy-CS filters out the lower quality regions, achieving high efficient bandwidth usage.

$\Theta(\cdot)$ is the feature encoder shared by all the collaborators and M represents the number of intermediate layers in $\Theta(\cdot)$.

Note that because each agent has its unique pose ξ^t , we need conduct feature alignment operation $\Lambda(\cdot, \cdot)$ across agents for subsequent collaboration process, *e.g.*, for the i -th agent, the transformed feature map from the k -th agent of the j -th intermediate layer is $F_{k \rightarrow i,j}^{e,t} \leftarrow \Lambda(F_{k,j}^{e,t}, \xi_{k \rightarrow i}^t)$, where the $\xi_{k \rightarrow i}^t$ is based on two ego poses ξ_k^t and ξ_i^t . Now the $F_{k \rightarrow i,j}^{e,t}$ and $F_{i,j}^{e,t}$ are supported in the same coordinate system. The goal of this paper is to optimize the unified collaborative perception framework and each agent outputs the prediction of perception task: $\hat{Y}^t = \sum_{i=1}^N \{\hat{y}_i^t | \mathbf{H}(\Psi(\sum_{j=1}^M \Upsilon(F_{i,j}^{e,t}, \sum_{k \neq i}^N \Omega(F_{k \rightarrow i,j}^{e,t}))))\}$, where $\Omega(\cdot)$, $\Upsilon(\cdot, \cdot)$, $\Psi(\cdot)$ represent the entropy-CS, G-CGRU, and MGFE module, respectively. The $\mathbf{H}(\cdot)$ denotes the headers for different perception tasks.

3.2. Entropy-based Communication Selection

The intuition of entropy-CS is low computational complexity and high interpretability to filter out the lower-quality regions. Based on above considerations, we extend the traditional entropy theory[6] toward two-dimensional

communication selection. Hence, the entropy-CS is non-parameter and high efficient to reduce the heavy bandwidth burden brought by multi-resolution feature information.

To further compress bandwidth requirements, the ego agent i and collaborator will first compress their encoder features $\mathbf{F}_{i,j}^{e,t}, \mathbf{F}_{k \rightarrow i,j}^{e,t} \in \mathbb{R}^{K,K,C}$ along channels into the query matrix $\mathbf{M}_{i,j}^{e,t}, \mathbf{M}_{k \rightarrow i,j}^{e,t} \in \mathbb{R}^{K,K}$ through a trainable 1×1 convolution kernel $\mathbf{W}_{1 \times 1}$, which is supervised by the downstream loss function. Then, we construct the two-stage entropy-based selective-region communication mechanism:

Self-select stage We first let the k -th collaborator determine location $\mathbf{S}_{k \rightarrow k,j}^{e,t}$ itself, where potentially contains useful information for later collaboration:

$$\mathbf{S}_{k \rightarrow k,j}^{e,t} = \Gamma(\Phi(\mathbf{M}_{k \rightarrow i,j}^{e,t}, \mathbf{M}_{k \rightarrow i,j}^{e,t}), \delta_s) \quad (1)$$

where the $\Gamma(\cdot, \delta_s)$ is an element-wise function, zeroing out elements smaller than the values of top- $\delta_s\%$ in \cdot , and Φ is the entropy estimation function, which measures the correlation from $\mathbf{K} \in \mathbb{R}^{K,K}$ to $\mathbf{Q} \in \mathbb{R}^{K,K}$ at the location of \mathbf{L} . Note that if \mathbf{L} is not assigned, the Φ will compute all region of \mathbf{K}, \mathbf{Q} :

$$\Phi(\mathbf{K}, \mathbf{Q}, \mathbf{L}) = \{p^{mn} * \log(p^{mn})\}_{m,n \in \mathbf{L}}, \text{ with}$$

$$p^{mn} = \frac{1}{M * N} \sum_{j,k} \sigma(\mathbf{K}(m+j, n+k) - \mathbf{Q}(m, n))$$

where $\sigma(\cdot)$ is Sigmoid function. Note that if $\|\mathbf{S}_{k \rightarrow k,j}^{e,t}\| \simeq 0$, represents that the k -th collaborator almost has no sufficient information. The entropy-CS will close the later cross-select communication stage, which can further reduce the bandwidth requirements.

Cross-select stage After the self-entropy selection stage, we thus derive the cross-entropy selection to obtain the communication location $\mathbf{S}_{k \rightarrow i,j}^{e,t}$ with the broadcast $\mathbf{M}_{i,j}^{e,t}$:

$$\mathbf{S}_{k \rightarrow i,j}^{e,t} = \Gamma(\Phi(\mathbf{M}_{i,j}^{e,t}, \mathbf{M}_{k \rightarrow i,j}^{e,t}, \mathbf{S}_{k \rightarrow k,j}^{e,t}), \delta_c) \quad (2)$$

As demonstrated in Figure 3, the entries of the derived matrix $\mathbf{S}_{k \rightarrow i,j}^{e,t}$ indicate where to communicate, which requires about $\frac{1}{\delta_s \delta_c}$ bandwidth of existing SRAR-based methods[40, 19, 23]. However, the transmission $\mathbf{T}_{k \rightarrow i,j}^{e,t} \leftarrow \mathbf{F}_{k \rightarrow i,j}^{e,t}[\mathbf{S}_{k \rightarrow i,j}^{e,t}]$ is sparse, with many positions being 0. This could potentially harm further collaboration process[25]. To compensate for this, after the ego agent receives, we spatially interpolate these positions from their neighbors across all channels at local, using a similar approach in [43]: $\mathbf{F}_{k \rightarrow i,j}^{e,t} \leftarrow \text{Interpolate}(\mathbf{T}_{k \rightarrow i,j}^{e,t})$.

3.3. Graph-based Collaborative GRU

Once a requesting ego agent collects the information from its linked supporting collaborators, it is important to design how to integrate its local observation with the selected multi-resolution feature maps from supporters[19].

We observe that the detected agents in perception area

change slowly over time in the urban low-speed collaborative scenes[46, 17]. Hence, the collaborative feature maps from last time are actually helpful to the present process.

However, the general RNN model[35, 10] can only process one-dimensional features, which makes it hard to preserve the spatial structure and local details of the agent's observation.

Based on the above considerations, we propose a novel graph-based collaborative GRU module named G-CGRU for each resolution collaboration process, as demonstrated in Figure 2. The key intuition of G-CGRU is that the collaboration depends not only on supporters but also on requesting the ego agent's previous information.

For the ego agent i of the j -th resolution intermediate feature maps, the inputs of G-CGRU are hidden states $\mathbf{h}_{i,j}^{e,t-1} \in \mathbb{R}^{K,K,C}$, the ego agent observation $\mathbf{F}_{i,j}^{e,t}$, and the supporters' selected feature maps $\{\mathbf{F}_{k \rightarrow i,j}^{e,t}\}_{k \neq i}$. The module mainly has three key gates as follows:

Update and Reset gates The *Update* gate is responsible for determining where the previous information needs to be passed along the next state and the *Reset* is to decide how much of the past information is needed to neglect. The *Update* and *Reset* gates share the same structure but different weights. Here we take *Reset* as an examples:

$$\mathbf{W}_{ir} = \sigma(\mathbf{W}_{3 \times 3} * ([\mathbf{h}_{i,j}^{e,t-1}; \mathbf{F}_{i,j}^{e,t}]))$$

$$\mathbf{R}_{i,j}^t = \sigma(\mathbf{W}_{ir} \odot \hat{\mathbf{h}}_{i,j}^{e,t-1} + (1 - \mathbf{W}_{ir}) \odot \mathbf{F}_{i,j}^{e,t}) \quad (3)$$

where $\odot, \sigma(\cdot), [;\cdot]$ represent dot product, Sigmoid function, and concatenation operation along channel dimensions, respectively. $*$ indicate a 3×3 convolution operation. The $\mathbf{R}_{i,j}^t \in \mathbb{R}^{K,K,C}$ learns where the hidden features $\mathbf{h}_{i,j}^{e,t-1}$ are conducive to the present.

Collab gate Based on the above *Reset* and *Update* gates, we thus derive the *Collab* gate. To make better collaborative feature integration, we construct a collaboration graph $\mathcal{G}_c^t(\mathcal{V}, \mathcal{E})$ in *Collab* gate, where node $\mathcal{V} = \{\mathcal{V}_i\}_{i \in 1, \dots, N}$ is the set of collaborative agents with the real-time pose information in the environment and $\mathcal{E} = \{\mathbf{W}_{i \rightarrow j}\}_{i,j \in 1, \dots, N, i \neq j}$ is the set of trainable edge matrix weights between agents and models the collaboration strength between two agents. Let $\mathcal{C}_{\mathcal{G}_c^t}(\cdot)$ be the collaboration process defined in the *Collab* module's graph \mathcal{G}_c^t . The j -th resolution enhanced maps of ego i agent after collaboration are $\mathbf{E}_{i,j}^{e,t} \leftarrow \mathcal{C}_{\mathcal{G}_c^t}(\mathbf{h}_{i,j}^{e,t-1}, \mathbf{F}_{k \rightarrow i,j}^{e,t}, \mathbf{F}_{i,j}^{e,t})$. This process has two stages: message attention (**S1**) and message aggregation (**S2**).

In the *message attention stage* (**S1**), each agent determines the matrix-valued edge weights, which reflect the strength from one agent to another at each individual cell. To determine the edge weights, we firstly get the conducive history information from hidden features by *Reset* gates through $\hat{\mathbf{h}}_{i,j}^{e,t-1} \leftarrow \mathbf{h}_{i,j}^{e,t-1} \odot \mathbf{R}_{i,j}^t$. Then, we utilize the edge

encode Π to correlate the history information, the feature map from another agent and ego feature map; that is, the matrix-value edge weight from k -th agent to the i -th agent is $\mathbf{W}_{k \rightarrow i} = \Pi(\hat{\mathbf{h}}_{i,j}^{e,t-1}, \mathbf{F}'_{k \rightarrow i,j}, \mathbf{F}_{i,j}^{e,t}) \in \mathbb{R}^{K,K}$, where Π concatenates three feature maps along the channel dimension and then utilizes four 1×1 convolutional layers to gradually reduce the number of channels from $3C$ to 1. Also, to normalize the edge weights across different agents, we implement a softmax operation on each cell of the feature map. Compared with previous work[23, 12, 45] generally consider time-discrete edge weight to reflect the static collaboration strength between two agents; while we consider time-continuous edge weight $\mathbf{W}_{k \rightarrow i}$ which dynamically models the collaboration strength from the k -th agent to the i -th agent from $t - 1$ to t . In the *message aggregation stage* (S2), each agent aggregates the feature maps from collaborators based on the normalized matrix-valued edge weights, the updated feature map $\mathbf{C}_{i,j}^{e,t}$ is utilized by $\sum_{k=1}^N \mathbf{W}_{k \rightarrow i} \circ \mathbf{F}'_{k \rightarrow i,j}$, where \circ represents the dot production broadcasting along the channel dimension. Finally, the collaborative map is $\mathbf{E}_{i,j}^{e,t} = \mathbf{Z}_{i,j}^t \odot \mathbf{C}_{i,j}^{e,t} + (1 - \mathbf{Z}_{i,j}^t) \odot \mathbf{h}_{i,j}^{e,t-1}$. Note that the $\mathbf{Z}_{i,j}^t$ is generated by *Update* gate and \odot is the dot product.

3.4. Multi-grain Feature Enhancement Module

The final ego agent i collaborative feature maps $\{\mathbf{E}_{i,j}^{e,t}\}_{j=1,\dots,M}$ with different resolutions generated from G-CGRU contain various visual context information, and each of them can be used to yield the downstream task.

To make the best use of multi-grain feature maps, we further propose a multi-grain feature enhancement (MGFE) module, which utilizes coarse- to fine-grained collaborative feature maps to guide the reconstruction process from a local to global perspective, as demonstrated in Figure 4. The MGFE consists of two stages:

Stage one A 1×1 convolution is applied to adjust the feature space of the collaborative feature map $\mathbf{E}_{i,j}^{e,t}$ to the ego observed feature map $\mathbf{F}_{i,j}^{e,t}$, then through global pooling along channels to the adjusted collaborative feature to obtain highlighted informative regions, following multiplied with the local feature map. The output \mathbf{f}_s^1 as follows:

$$\mathbf{f}_s^1 = \mathcal{G}(\sigma(\mathbf{W}_{1 \times 1} * \mathbf{E}_{i,j}^{e,t} + \mathbf{b})) \circ \mathbf{F}_{i,j}^{e,t} \quad (4)$$

where \mathcal{G} , \circ denotes global max pooling operation along the channels and dot production broadcasting along the channel dimension, respectively. $\mathbf{W}_{1 \times 1}$ indicates trainable parameters, \mathbf{b} refers bias.

Stage two The coarse-grained, ego-observed feature map $\mathbf{F}_{i,j-1}^{e,t}$ is upsampled to the same dimension as the high-resolution map, and then concatenated with the \mathbf{f}_s^1 , $\mathbf{E}_{i,j}^{e,t}$. The output of stage two \mathbf{f}_s^2 as follows:

$$\mathbf{f}_s^2 = [L_2(\text{upsample}(\mathbf{F}_{i,j-1}^{e,t})); L_2(\mathbf{f}_s^1); L_2(\mathbf{E}_{i,j}^{e,t})] \quad (5)$$

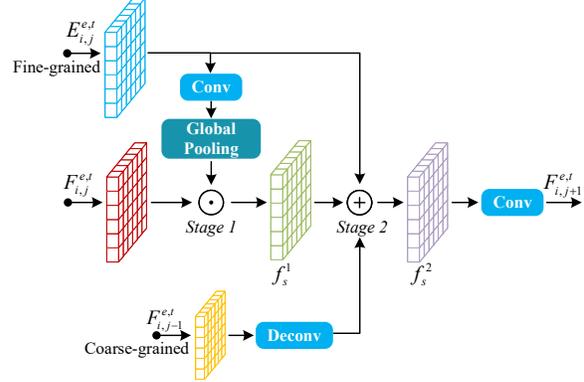


Figure 4. The architecture of multi-grain feature enhancement module. The fine- and coarse-grained feature maps guide the reconstruction process from the local to global perspective, respectively.

where *upsample* denotes the deconvolution operation, L_2 refers to the L_2 norm layer[22], which is helpful for combining two (or more) feature maps. To adjust the channels dimension of \mathbf{f}_s^2 , a 3×3 convolution is applied to \mathbf{f}_s^2 : $\mathbf{F}_{i,j+1}^{e,t} \leftarrow \sigma(\mathbf{W}_{3 \times 3} * \mathbf{f}_s^2 + \mathbf{b})$.

After the reconstruction process of the MGFE module, we can obtain the ego agent enhancement feature $\mathbf{D}_i^{e,t}$. To obtain the final detection outputs $\tilde{\mathbf{Y}}_i^t$, we employ an output header $\mathbf{H}(\cdot)$: $\tilde{\mathbf{Y}}_i^t \leftarrow \mathbf{H}(\mathbf{D}_i^{e,t})$. To implement the header $\mathbf{H}(\cdot)$, we employ two branches of convolutional layers to classify the foreground-background categories and regress the bounding boxes.

To train our model, we employ the binary cross-entropy loss to supervise foreground-background classification and the smooth L_1 loss to supervise the bounding-box regression:

$$\mathcal{L} = \sum_{i=1}^N \mathcal{L}_{det}(\mathbf{Y}_i, \tilde{\mathbf{Y}}_i) \quad (6)$$

where \mathcal{L}_{det} denotes the classification and regression losses and $\mathbf{Y}_i, \tilde{\mathbf{Y}}_i$ represents the ground-truth detection and predicted detection in the perception region of the i -th agent, respectively.

4. Experiment

4.1. Collaborative 3D Object detection dataset

To validate our proposed UMC on the LIDAR-based 3D object detection task. Same as [19, 15, 20], we utilize the multi-agent datasets of V2X-Sim[17] and OPV2V[46]. The V2X-Sim is built by the co-simulation of SUMO[26] and CARLA[8], which contains 90 scenes, and each scene includes 100 frames. The OPV2V is co-simulated by OpenCDA[46] and CARLA, including 12K frames of 3D point clouds with 230K annotated 3D boxes.

Table 1. Detection comparison on V2X-Sim dataset. Key: [Best, Second Best, Third Best]

Models	ARSV _{50/70} ↑	ARCV _{50/70} ↑	ARCI _{50/70} ↓	ARTC _{50/70} ↑	AP _{50/70} ↑
No Fusion	76.90/72.35	4.62/3.45	2.30/1.27	7.85/4.62	51.87/45.05
Early Fusion	86.70/83.72	64.17/60.49	9.55/7.88	15.67/13.33	67.79/62.29
Who2com	74.60/69.30	4.70/3.95	1.80/1.30	7.58/7.25	49.77/42.30
When2com	75.10/69.69	4.79/4.08	1.83/1.28	7.58/7.25	52.06/44.21
V2VNet	80.13/75.17	44.66/32.00	1.33/0.97	10.70/6.83	58.51/48.98
DiscoNet	86.64/82.16	58.58/46.47	3.02/1.91	12.86/9.32	65.89/56.74
Where2comm	81.96/77.24	44.97/30.90	11.45/8.45	16.51/11.61	62.05/54.59
UMC(ours)	84.67/80.68	67.01/60.04	2.38/1.70	12.35/9.46	67.80/60.01

Table 2. Detection comparison on OPV2V dataset. Key: [Best, Second Best, Third Best]

Models	ARSV _{50/70} ↑	ARCV _{50/70} ↑	AP _{50/70} ↑
No Fusion	69.33/46.35	12.32/4.36	54.69/23.94
Early Fusion	64.62/46.95	45.00/24.39	55.88/25.89
Who2com	69.43/45.68	15.84/6.21	55.36/23.39
When2com	69.27/45.54	15.90/6.22	55.13/23.27
V2VNet	73.42/48.21	24.65/11.16	59.03/24.71
DiscoNet	72.04/46.94	37.25/22.57	55.46/23.04
Where2comm	69.79/47.60	13.89/6.09	56.02/25.38
UMC(ours)	76.56/47.68	47.82/25.06	61.90/24.50

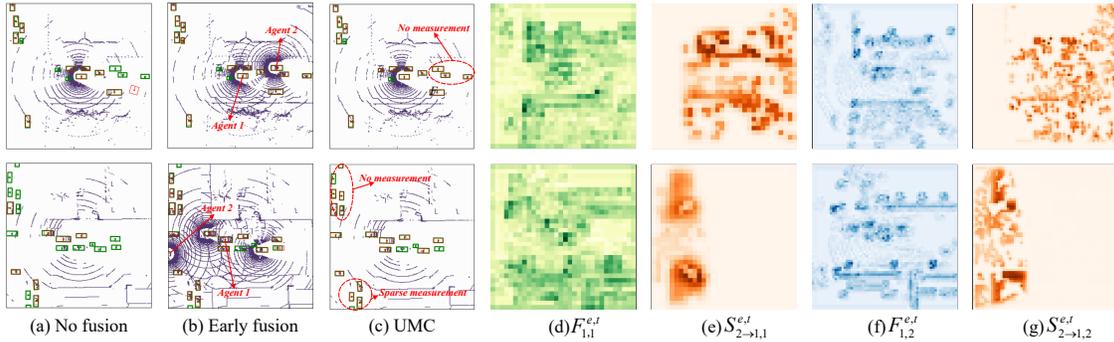


Figure 5. Detection and communication selection for Agent 1. The green and red boxes represent the ground truth (GT) and predictions, respectively. (a-c) shows the results of no fusion, early fusion, and UMC compared to GT. (d) The coarse-grained collaborative feature of Agent 1. (e) Matrix-valued entropy-based selected communication coarse-grained feature map from Agent 2. (f) The fine-grained collaborative feature of Agent 1. (g) Matrix-valued entropy-based selected communication fine-grained feature map from Agent 2.

4.2. Implementation details

Experimental setting Same as [19, 15, 46], we crop the point clouds located in the region of $[-32m, 32m] \times [-32m, 32m] \times [0, 5m]$ defined in the ego-vehicle Cartesian coordinate system. The size of each voxel cell is set as $0.25m \times 0.25m \times 0.4m$ to get the Bird’s-Eyes view map with dimension $256 \times 256 \times 13$. The SRAR-based transmitted collaborative feature map (TCF) has a dimension of $32 \times 32 \times 256$. Our proposed UMC’s TCFs have the dimension of $32 \times 32 \times 256, 64 \times 64 \times 128$. We train all the models using NVIDIA RTX 3090 GPU with PyTorch[27].

Baselines We consider the single-agent perception model, called *No Fusion*, which only processes a single-view point cloud. Also, we consider the holistic-view perception model based on early collaboration, called *Early Fusion*. We consider five collaboration methods: Who2com[24], When2com[23], Where2comm[12], DiscoNet[19] and V2VNet[40]. To ensure fairness, all the methods share the same encoder and detection header architecture and collaborate at the same intermediate feature layer.

Evaluation Metrics We employ the Average Precision (AP) to validate the overall performance of the model. To systematically analyze the performance of the model, we introduce new metrics from four aspects: i) **ARSV** denotes the Average **R**ecall of agents that are visible from **S**inge-vehicle **V**iew; ii) **ARCV** denotes the Average **R**ecall of agents that are invisible from single-vehicle view but visible from **C**ollaborative-**V**iew, which evaluates model’s

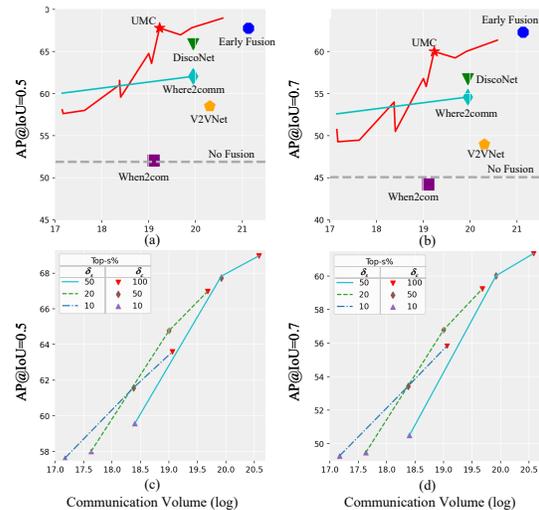


Figure 6. Performance-bandwidth trade-off in log scale.

performance of collaborative process; iii) **ARCI** denotes the Average **R**ecall of **C**ompletely-**I**nvisible agents, which evaluates the performance of detecting agents with limited information; iv) **ARTC** denotes the Average **R**ecall of the agents that were visible at the last time but not visible at this time, which evaluates model’s performance of **T**ime **C**ontinuity. Note that more technical details are shown in Appendix.1, and we employ all evaluation metrics at the Intersection-over-Union (IoU) threshold of 0.5 and 0.7.

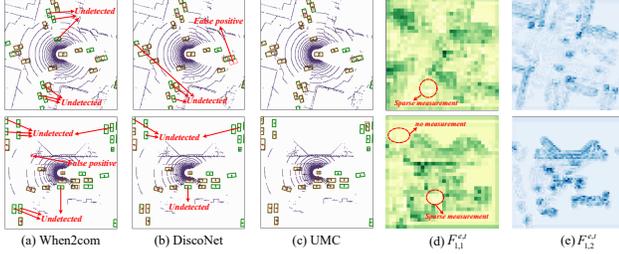


Figure 7. UMC qualitatively outperforms the state-of-the-art methods. The green and red boxes denote ground truth and detection, respectively. (a) Results of When2com. (b) Results of DiscoNet. (c) Results of UMC. (d)-(e) Agent 1’s coarse-grained and fine-grained collaborative feature maps, respectively.

4.3. Quantitative evaluation

Overall performance comparison Table 1 of $AP_{50/70}$ column shows the comparisons in terms of $AP(@IoU=0.5/0.7)$. We see that i) early fusion achieves the best detection performance when $AP@0.7$, and there is an obvious improvement over no fusion, *i.e.*, $AP@0.5$ and $AP@0.7$ are increased by 30.69% and 38.26% respectively, revealing the effectiveness of collaboration; ii) among the SRAR-based intermediate collaboration methods, the proposed UMC achieves the best performance. Compared to When2com, UMC improves by 30.23% in $AP@0.5$ and 35.71% in $AP@0.7$. Compared to the knowledge distillation (KD) based DiscoNet, UMC improves by 2.9% in $AP@0.5$ and 5.76% in $AP@0.7$.

Proposed new metric performance comparison From the rest column of Tabel 1, we see that i) the ARSV reflects that collaboration can improve the performance of original single-view based objects. The early fusion achieves the best performance, and there is an obvious improvement over no fusion, *i.e.*, $ARSV@0.5$ and $ARSV@0.7$ are increased by 12.74% and 15.72% respectively. Note that the DiscoNet inherits the performance of early fusion through knowledge distillation; ii) the ARCV reflects the quality of the collaboration process. The proposed UMC achieves the best performance. Compared to the DiscoNet, UMC improves by 14.39% in $ARCV@0.5$ and 29.07% in $ARCV@0.7$. Note that in when2com and who2com, the collaboration is simply integrated by stacking features along channels, resulting in the low quality of collaboration; iii) the ARCI reflects the detection performance under limited information. Thus, a high ARCI may cause more false positives, reducing the precision of the model; iv) the ARTC reflects the time continuity of the model, and the ARTC of existing methods is relatively low, In that sense, our work may unlock the future temporal reasoning in collaborative perception.

Performance-bandwidth trade-off analysis As demonstrated in Figure 6, we comprehensively compare the proposed UMC under different (δ_s, δ_c) values with the baseline

Table 3. Detection performance with different grains selection. $F_{i,m}^{e,t}$ ($n = 1, 2, 3$) denotes different grained collaborative features. Note that the C. V. is the short for Communication Volume.

Multi-Grains Selection			$AP_{50/70} \uparrow$	C. V. \downarrow
$F_{i,1}^{e,t}$	$F_{i,2}^{e,t}$	$F_{i,3}^{e,t}$		
✓	✓		67.80/60.01	19.23
✓		✓	58.36/50.71	19.84
	✓	✓	60.27/53.63	20.02

methods in terms of the trade-off between performance and communication bandwidth. The no fusion model is shown as a dashed line since the communication volume is zero. We see that i) UMC achieves the best trade-off among the other methods; ii) the detect performance of where2comm is more stable than UMC when bandwidth is reduced. The where2comm is specific for the detection downstream task, which is based on the sparsity of foreground information to reduce communication by filtering the background intensively with a small performance sacrifice. And, our entropy-CS aims to optimize not only detection but also general downstream tasks based on the traditional information theory. Hence, when setting small value of (δ_s, δ_c) , the detection performance will degrade more easily than where2comm, shown in Figure 6 (c)-(d).

4.4. Qualitative evaluation

Visualization of UMC working mechanism To understand the working mechanism of the proposed UMC, we visualize the detection results, different grains collaborative feature maps, and the corresponding entropy-based transmission maps; see Figure 5. Note that the proposed entropy-based selection is matrix-based, reflecting the informative regions in a cell-level resolution, which is shown as the 2D communication map that is compressed along channels by a sum operation. Figure 5 (a), (b), and (c) represent three detection results of Agent 1 based on the no fusion, early fusion, and the proposed UMC, respectively. We see that with collaboration, UMC is able to detect objects in these sheltered and long-range zones. To further explain the rationale behind the results, Figure 5 (d) and (f) provide the corresponding ego coarse- and fine-grained collaborative feature maps, respectively. We clearly observe that the coarse-grained feature reflects the global structure information and the fine-grained feature contains rich local texture details, the proposed UMC complements their advantages. Figure 5 (e) and (g) show the coarse- and fine-grained communication maps from agent 2 to agent 1, respectively. We observe that with the grains being finer, the communication map becomes more sparse, meaning the entropy-based selection can more precisely find the informative regions relative to ego agents, also, the entropy-based selection ensures lower bandwidth under the multi-grain collaborative process.

Comparison with DiscoNet and When2com Figure 7 shows two examples to compare the proposed UMC with

Table 4. Quantitative results of ablation experiments on V2X-Sim.

Variants	Entropy-CS	G-CGRU	MGFE	ARSV _{50/70} ↑	ARCV _{50/70} ↑	ARCI _{50/70} ↓	ARTC _{50/70} ↑	AP _{50/70} ↑
(1)	✓	✓	✓	86.67/80.68	67.01/60.04	2.38/1.70	12.35/9.46	67.80/60.01
(2)		✓	✓	87.67/83.29	67.88/62.17	2.87/1.66	12.01/9.55	65.72/58.78
(3)	✓	✓		79.38/74.90	40.35/34.14	1.49/1.28	3.79/3.52	56.73/49.22
(4)		✓		87.72/82.13	60.58/46.53	3.13/1.43	12.78/9.98	67.30/56.19

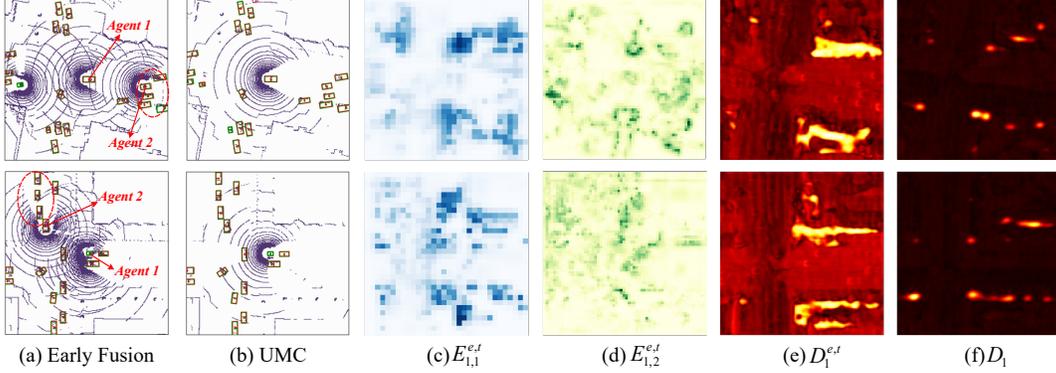


Figure 8. Detection and feature map visualization for Agent 1. Green/red boxes represent GT/predictions. (a) and (b) denote the results of early fusion and UMC. (c) and (d) denote the coarse- and fine-grained collaborative maps generated from the G-CGRU module, respectively. (e) and (f) denote the UMC and early fusion’s downstream features D_1 for final detection output \hat{Y}_1 , respectively.

When2com and DiscoNet. We see that UMC is able to detect more objects, especially in edge regions with sparse measurement. The reason is that both When2com and DiscoNet employ only a single coarse-grained global feature map for collaborations, which loses many details of the scene, while UMC combines coarse- and fine-grained features for collaboration, see the visualization in Figure 7 (d)-(e). Remarkably, compared to knowledge-distillation-based DiscoNet, the UMC does not need to pretrain different teacher models (early fusion) when switching to different datasets, and the performance will be influenced greatly when the teacher model is not performing well, as shown in Table 2.

4.5. Ablation study

We validate several key designs of our framework, including the entropy-based communication selection, the G-CGRU, and the MGFE module. The quantitative results are shown in Table 4.

Effect of grains selection In table 3 we investigate the detection performance when employing different grain selections. Note that the shape of $F_{i,n}^{e,t} (n = 1, 2, 3) \in \mathbb{R}^{W,H,C}$ is $32 \times 32 \times 256$, $64 \times 64 \times 128$ and $128 \times 128 \times 64$, respectively. We see that i) applying collaborative feature maps on the selection of $(F_{i,1}^{e,t}, F_{i,2}^{e,t})$ has the best trade-off between performance and bandwidth; and ii) the selection of $(F_{i,1}^{e,t}, F_{i,3}^{e,t})$ for collaboration has the worst detection performance. This is because the grains of feature difference between the two is too large, resulting in deviation in feature fusion process; and iii) the performance of selection of $(F_{i,2}^{e,t}, F_{i,3}^{e,t})$ is slightly worse than $(F_{i,1}^{e,t}, F_{i,2}^{e,t})$, because

both $(F_{i,2}^{e,t}, F_{i,3}^{e,t})$ focus on the local texture details, resulting in the loss of global structure information; and iv) on the premise of performance of $(F_{i,1}^{e,t}, F_{i,2}^{e,t})$, the heavy bandwidth of the $(F_{i,1}^{e,t}, F_{i,2}^{e,t}, F_{i,3}^{e,t})$ selection makes the potential performance improvement meaningless.

Effect of collaboration strategy Table 4 compares the different ablated networks. We see that i) the UMC (variants 1) achieves the best detection performance among the other ablated networks; and ii) variant 4 reflects that the redesigned graph-based collaborative GRU module has an obvious improvement over no fusion, *i.e.*, AP@0.5 and AP@0.7 are increased by 26.66% and 24.72% respectively, see the Figure 8 (c) and (d); and iii) variant 2 is composed of variants 4 and multi-grain feature enhance module. Compared to variant 4, variant 2 improves by 12.05% in ARCV@0.5 and 33.61% in ARCV@0.7, showing the effectiveness of the MGFE module. Meanwhile, as demonstrated in Figure 8 (e) and (f), the MGFE module can enhance the downstream feature map compared to the early fusion method; and iv) from the results of variants 3 and 4, and the results of variants 1 and 2, we observe that the MGFE module can give positive guidance to the trainable entropy-based communication selection module because the MGFE module allows the entropy-CS to comprehensively select regions from local and global viewpoints.

5. Conclusion

We propose a unified, bandwidth-efficient collaborative perception framework named UMC. Its core concept is to utilize the multi-resolution technique to cover the intrinsic

interactions among the communication, collaboration, and reconstruction processes in MCP. To validate it, we conduct experiments on the V2X-Sim and OPV2V datasets and propose a brand new evaluation metric. Comprehensive quantitative and qualitative experiments show that the UMC achieves an outstanding performance-bandwidth trade-off among the existing collaborative perception methods.

References

- [1] Qi Chen. F-cooper: feature based cooperative perception for autonomous vehicle edge computing system using 3d point clouds. *Proceedings of the 4th ACM/IEEE Symposium on Edge Computing*, 2019.
- [2] Qi Chen, Sihai Tang, Qing Yang, and Song Fu. Cooper: Cooperative perception for connected autonomous vehicles based on 3d point clouds. *2019 IEEE 39th International Conference on Distributed Computing Systems (ICDCS)*, pages 514–524, 2019.
- [3] Shanzhi Chen, Jinling Hu, Yan Shi, Li Zhao, and Wen Li. A vision of c-v2x: Technologies, field testing, and challenges with chinese development. *IEEE Internet of Things Journal*, 7(5):3872–3881, 2020.
- [4] Wei-Ting Chen, H. Fang, Cheng-Lin Hsieh, Cheng-Che Tsai, I-Hsiang Chen, Jianwei Ding, and Sy-Yen Kuo. All snow removed: Single image desnowing algorithm using hierarchical dual-tree complex wavelet representation and contradict channel loss. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4176–4185, 2021.
- [5] Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder–decoder approaches. In *SSST@EMNLP*, 2014.
- [6] Gavin E Crooks. On measures of entropy and information. *Tech. Note*, 9:v4, 2017.
- [7] Jiaxun Cui, Hang Qiu, Dian Chen, Peter Stone, and Yuke Zhu. Coopernaut: End-to-end driving with cooperative perception for networked vehicles. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [8] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. CARLA: An open urban driving simulator. In *Proceedings of the 1st Annual Conference on Robot Learning*, pages 1–16, 2017.
- [9] Lingxiao He, Jian Liang, Haiqing Li, and Zhenan Sun. Deep spatial feature reconstruction for partial person re-identification: Alignment-free approach. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7073–7082, 2018.
- [10] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [11] Xin Hong, Pengfei Xiong, Renhe Ji, and Haoqiang Fan. Deep fusion network for image completion. In *Proceedings of the 27th ACM International Conference on Multimedia*, MM '19, page 2033–2042, New York, NY, USA, 2019. Association for Computing Machinery.
- [12] Yue Hu, Shaoheng Fang, Zixing Lei, Yiqi Zhong, and Siheng Chen. Where2comm: Communication-efficient collaborative perception via spatial confidence maps. In *Thirty-sixth Conference on Neural Information Processing Systems (Neurips)*, November 2022.
- [13] Hongwu Kuang, Bei Wang, Jianping An, Ming Zhang, and Zehan Zhang. Voxel-fpn: Multi-scale voxel feature aggregation for 3d object detection from lidar point clouds. *Sensors (Basel, Switzerland)*, 20, 2020.
- [14] Alex H. Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12689–12697, 2019.
- [15] Zixing Lei, Shunli Ren, Yue Hu, Wenjun Zhang, and Siheng Chen. Latency-aware collaborative perception. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022.
- [16] Ruoteng Li, Robby T. Tan, and Loong Fah Cheong. All in one bad weather removal using architectural search. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3172–3182, 2020.
- [17] Yiming Li, Dekun Ma, Ziyang An, Zixun Wang, Yiqi Zhong, Siheng Chen, and Chen Feng. V2x-sim: Multi-agent collaborative perception dataset and benchmark for autonomous driving. *IEEE Robotics and Automation Letters*, 7(4):10914–10921, 2022.
- [18] Yiming Li, Dekun Ma, Ziyang An, Zixun Wang, Yiqi Zhong, Siheng Chen, and Chen Feng. V2x-sim: Multi-agent collaborative perception dataset and benchmark for autonomous driving. *IEEE Robotics and Automation Letters*, 7:10914–10921, 2022.
- [19] Yiming Li, Shunli Ren, Pengxiang Wu, Siheng Chen, Chen Feng, and Wenjun Zhang. Learning distilled collaboration graph for multi-agent perception. In *NeurIPS*, 2021.
- [20] Yiming Li, Juexiao Zhang, Dekun Ma, Yue Wang, and Chen Feng. Multi-robot scene completion: Towards task-agnostic collaborative perception. In *6th Annual Conference on Robot Learning*, 2022.
- [21] Tsung-Yi Lin, Piotr Dollár, Ross B. Girshick, Kaiming He, Bharath Hariharan, and Serge J. Belongie. Feature pyramid networks for object detection. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 936–944, 2017.
- [22] Wei Liu, Andrew Rabinovich, and Alexander C. Berg. Parsenet: Looking wider to see better. *ArXiv*, abs/1506.04579, 2015.
- [23] Yen-Cheng Liu, Junjiao Tian, Nathan Glaser, and Zsolt Kira. When2com: Multi-agent perception via communication graph grouping. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4105–4114, 2020.
- [24] Yen-Cheng Liu, Junjiao Tian, Chih-Yao Ma, Nathan Glaser, Chia-Wen Kuo, and Zsolt Kira. Who2com: Collaborative perception via learnable handshake communication. *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6876–6883, 2020.
- [25] Zhuang Liu, Zhiqiu Xu, Hung-Ju Wang, Trevor Darrell, and Evan Shelhamer. Anytime dense prediction with confidence

- adaptivity. *International Conference on Learning Representations (ICLR)*, 2022.
- [26] Pablo Alvarez Lopez, Michael Behrisch, Laura Bieker-Walz, Jakob Erdmann, Yun-Pang Flötteröd, Robert Hilbrich, Leonhard Lücken, Johannes Rummel, Peter Wagner, and Evarmarie Wießner. Microscopic traffic simulation using sumo. In *The 21st IEEE International Conference on Intelligent Transportation Systems*. IEEE, 2018.
- [27] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017.
- [28] Jialun Peng, Dong Liu, Songcen Xu, and Houqiang Li. Generating diverse structure for image inpainting with hierarchical vq-vae. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10770–10779, 2021.
- [29] C. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 77–85, 2017.
- [30] C. Qi, L. Yi, Hao Su, and Leonidas J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *NIPS*, 2017.
- [31] Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 779–788, 2016.
- [32] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39:1137–1149, 2015.
- [33] 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*, pages 234–241. Springer, 2015.
- [34] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. *ArXiv*, abs/1505.04597, 2015.
- [35] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- [36] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. Pv-rcnn: Point-voxel feature set abstraction for 3d object detection. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10526–10535, 2020.
- [37] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointcnn: 3d object proposal generation and detection from point cloud. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–779, 2019.
- [38] Martin Simon, Stefan Milz, Karl Amende, and Horst-Michael Groß. Complex-yolo: An euler-region-proposal for real-time 3d object detection on point clouds. In *ECCV Workshops*, 2018.
- [39] Angtian Wang, Yihong Sun, Adam Kortylewski, and Alan Loddon Yuille. Robust object detection under occlusion with context-aware compositionalsnets. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12642–12651, 2020.
- [40] Tsun-Hsuan Wang, Sivabalan Manivasagam, Ming Liang, Binh Yang, Wenyuan Zeng, James Tu, and Raquel Urtasun. V2vnet: Vehicle-to-vehicle communication for joint perception and prediction. In *ECCV*, 2020.
- [41] Wentao Wang, Jianfu Zhang, Li Niu, Haoyu Ling, Xue Yang, and Liqing Zhang. Parallel multi-resolution fusion network for image inpainting. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14539–14548, 2021.
- [42] Pengxiang Wu and Siheng Chen. Motionnet: Joint perception and motion prediction for autonomous driving based on bird’s eye view maps. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11382–11392, June 2020.
- [43] Zhenda Xie, Zheng Zhang, Xizhou Zhu, Gao Huang, and Steve Lin. Spatially adaptive inference with stochastic feature sampling and interpolation. In *European Conference on Computer Vision (ECCV)*, August 2020.
- [44] Runsheng Xu, Zhengzhong Tu, Hao Xiang, Wei Shao, Bolei Zhou, and Jiaqi Ma. Cobevt: Cooperative bird’s eye view semantic segmentation with sparse transformers. In *Conference on Robot Learning (CoRL)*, 2022.
- [45] Runsheng Xu, Hao Xiang, Zhengzhong Tu, Xin Xia, Ming-Hsuan Yang, and Jiaqi Ma. V2x-vit: Vehicle-to-everything cooperative perception with vision transformer. *ArXiv*, abs/2203.10638, 2022.
- [46] Runsheng Xu, Hao Xiang, Xin Xia, Xu Han, Jinlong Li, and Jiaqi Ma. Opv2v: An open benchmark dataset and fusion pipeline for perception with vehicle-to-vehicle communication. In *ICRA*, 2022.
- [47] Yan Yan, Yuxing Mao, and Bo Li. Second: Sparsely embedded convolutional detection. *Sensors (Basel, Switzerland)*, 18, 2018.
- [48] Binh Yang, Wenjie Luo, and Raquel Urtasun. Pixor: Real-time 3d object detection from point clouds. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7652–7660, 2018.
- [49] Haibao Yu, Yizhen Luo, Mao Shu, Yiyi Huo, Zebang Yang, Yifeng Shi, Zhenglong Guo, Hanyu Li, Xing Hu, Jirui Yuan, and Zaiqing Nie. Dair-v2x: A large-scale dataset for vehicle-infrastructure cooperative 3d object detection. In *IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, June 2022.
- [50] Xiaoding Yuan, Adam Kortylewski, Yihong Sun, and Alan Loddon Yuille. Robust instance segmentation through reasoning about multi-object occlusion. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11136–11145, 2021.
- [51] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Multi-stage progressive image restoration. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14816–14826, 2021.

- [52] Yanhong Zeng, Jianlong Fu, Hongyang Chao, and Baining Guo. Learning pyramid-context encoder network for high-quality image inpainting. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1486–1494, 2019.
- [53] Kaihao Zhang, Rongqing Li, Yanjiang Yu, Wenhan Luo, and Changsheng Li. Deep dense multi-scale network for snow removal using semantic and depth priors. *IEEE Transactions on Image Processing*, 30:7419–7431, 2021.
- [54] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4490–4499, 2018.

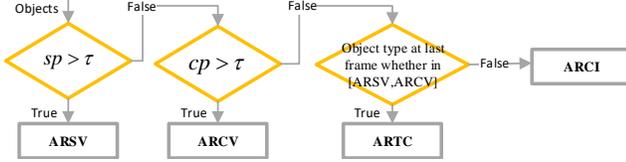


Figure 9. In the test set, we traverse all objects in each frame to obtain their corresponding types. Note that sp , cp are short for detected points from single and collaborative view, respectively.

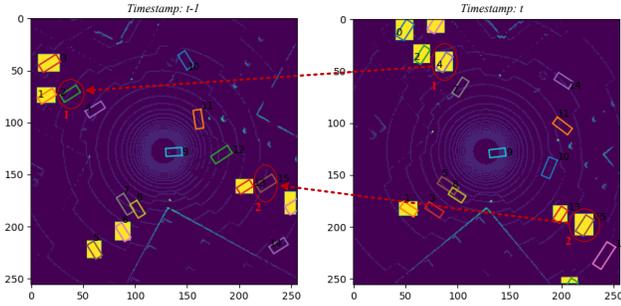


Figure 10. The detailed process of manually labelling.

A. Detailed information of Proposed Metrics

The detailed process of proposed metrics are shown in Figure 9. We traverse all objects in each frame to obtain their corresponding types. Note that the type of ARSV/ARCV can be automatically generated with code, as shown in Listing 1. And the type of ARCI/ARTC needs to be manually labeled, because whether it is visible at the last moment cannot be directly determined due to the movement of the vehicles from $t - 1$ to t , as shown in Figure 10.

For detailed process of manually labelling, *e.g.*, at timestamp t , we firstly label vehicles with yellow background as ARCI type. Then, by comparing with time $t - 1$, we find that agent 1,2 are visible at time $t - 1$ (with no yellow background). Based on that, we change the type of agent 1,2 from ARCI to ARTC. The pseudo code are as follows:

```

1 single_view = np.load(...)
2 collab_view = np.load(...)
3 for object in len(single_view.objects):
4     if single_view[object].points > threshold:
5         object.type = ARSV
6     else:
7         if collab_view[object].points > threshold:
8             object.type = ARCV
9         else:
10            object.type = ARCI
11
12 % Finally, we will manually label partial ARCI
    to ARTC, as shown in Figure 2.

```

Listing 1. Pseudo code for proposed metrics

As for the threshold points τ , the proposed metrics distinguish between ARSV and ARCV based on whether they are visible. Nevertheless, whether they are visible depends on the number of points included in each object and the per-

Table 5. The ARSV and ARCV performance with different threshold points τ .

Points	ARSV _{50/70}	ARCV _{50/70}
10	81.48/77.48	23.72/19.34
7	79.44/75.05	15.34/11.88
5	77.73/73.21	6.08/4.07
4	76.90/72.35	4.62/3.45

formance of the detector. To verify at how many points the detector can not detect the object, we conducted the following experiments on the *No Fusion* model.

Table 5 shows the ARSV and ARCV in terms of different points. We can see that i) as the number of points decreases, so does the ARCV. This is because, as the points become more accurate, the no fusion model should theoretically be 0 in terms of ARCV; ii) when the points are greater than 5, the decline in ARCV is very large. When points are less than 5, the decline of the ARCV slows down, indicating it is close to the accurate points; iii) when points equal 4, the ARCV is 4.62% in IoU@0.5 and 3.45% in IoU@0.7, which are within the acceptable error range of 5%. So, to decide if an object is visible, we look at how many points it has and whether that number is greater than 4.

Note that the ‘last frame’ of ARTC’s time interval varies based on the sampling frequency (5Hz in V2X-Sim, 10Hz in OPV2V). And, the ARSV and ARCV only take recall into account, while the AP is weighted by both recall and precision. Hence, a high AP does not necessarily mean high ARSV or ARCV values. Based on that, you may wonder why not utilize APSV/APCV as the new metric. It is because the each proposed model can only predict the classifications (background or foreground) and regressions (x,y,w,h) of each pixel. Therefore, there is no prediction of the corresponding types for each detected objects.

B. Experiments details

B.1. Basic parameters

Our experiments are all performed on the workstation with AMD Core Ryzen Threadripper 3960X CPU and Nvidia 3090 GPU with Pytorch v1.7.1, CUDA 11.0. The SRAR-based’s transmitted collaborative feature map (TCF) has a dimension of $32 \times 32 \times 256$. Our proposed UMC’s TCFs have the dimension of $32 \times 32 \times 256$, $64 \times 64 \times 128$. As for hyper-parameter tuning, we choose Adam as the optimizer and set the batch size to 4 for both V2X-Sim and OPV2V datasets. Also, we utilize the same number of scenes (total 80 scenes) for training. For testing stage, the V2X-sim utilizes 10 scenes, and OPV2V utilizes 15 scenes. Meanwhile, we use initial learning rate of 0.001 and set the random seed to 622.

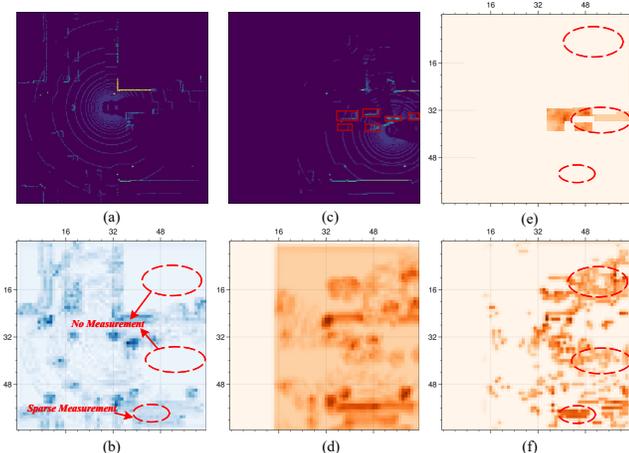


Figure 11. Difference between where2comm and UMC communication strategy. (a) The ego agent’s observation. (b) The ego agent’s observation at feature level. (c) The collaborator’s observation, red box denoted for the detected agents. (d) The collaborator’s observation at feature level. (e) Communication map of where2comm. (f) Communication map of UMC.

B.2. Baseline setting

To ensure fairness, we fix the structure of the shared feature extractor MotionNet[42] and detector, and transplant the collaborative part of the different methods without modification. Meanwhile, all the models are trained 100 epoch with initial learning rate of 0.001 and set the learning rate update strategy as ‘torch.optim.lr_scheduler.MultiStepLR(self. optimizer_head, milestones=[50, 100], gamma=0.5)’.

B.3. Communication Volume

The communication volume for each method is calculated by: $\text{mean}(\log(\sum_{i=1}^{\text{agents}} \sum_{t=1}^T (F.\text{size} + Q.\text{size})))$ during the test stage. F represents the transmitted feature map, and Q represents the query matrix, which is calculated in UMC, who2com, and when2com, and equals 0 in other methods.

B.4. Setting of δ_s, δ_c

Table 6. The performance-bandwidth trade off with different δ_s, δ_c values on V2X-Sim dataset.

δ_s	δ_c	ARSV _{50/70} ↑	ARCV _{50/70} ↑	ARCI _{50/70} ↓	ARTC _{50/70} ↑	AP _{50/70} ↑	C.V. ↓
50	100	85.66/81.87	70.76/63.15	2.41/1.7	11.88/8.12	68.97/61.35	20.58
50	50	84.78/80.73	67.46/60.72	2.50/1.78	11.88/8.12	67.83/60.02	19.92
50	10	80.62/75.39	31.48/23.53	2.17/1.37	12.03/7.20	59.57/50.50	18.4
20	100	87.30/84.08	74.08/66.77	16.68/12.24	22.15/17.41	66.98/59.24	19.68
20	50	83.23/78.98	56.49/48.87	2.12/1.51	11.23/8.08	64.77/56.79	19
20	10	80.53/75.38	19.02/12.33	2.14/1.38	9.16/6.74	58.00/49.46	17.63
10	100	86.57/82.97	59.37/51.68	14.57/10.51	19.09/15.15	63.59/55.84	19.06
10	50	82.12/77.44	38.52/31.28	2.16/1.47	10.77/8.35	61.57/53.43	18.376
10	10	80.78/75.51	13.29/8.97	2.21/1.39	10.77/8.35	57.62/49.27	17.171
5	100	86.04/82.59	47.19/40.44	14.17/9.9	19.25/15.14	61.04/53.99	18.375
5	20	80.87/75.72	15.09/10.76	2.12/1.36	10.77/8.35	57.67/49.45	17.17
1	100	85.44/81.49	34.59/28.72	13.47/9.88	20.49/17.35	58.07/50.72	17.15

We record the proposed UMC under different (δ_s, δ_c) in

terms of the trade-off between performance and communication bandwidth, as shown in Table 6.

Meanwhile, we comprehensively analyze the communication strategy between where2comm[12] and our proposed UMC. As shown in Figure 11, where2comm takes the advantages of sparsity of foreground information and only transmits the regions that the agents have. However, as for different downstream tasks, such as segmentation[44] or scene completion[20], there are other no-measurement or sparse-measurement regions that need collaborator’s communication, as shown in Figure 11.(b). Hence, our proposed UMC aim to optimize not only detection but also general downstream tasks based on the traditional information theory. From the Figure 11.(f), we can observe that UMC can transmit the necessary regions to ego agent for general downstream tasks.

Note that in addition to discussing the top- $\delta\%$ based filtering strategy of Eq.1 in manuscript, we also explored the mean based filtering strategy, the corresponding performance is shown as follows:

Table 7. The performance of mean based filtering strategy on V2X-Sim dataset.

δ_s	δ_c	ARSV _{50/70} ↑	ARCV _{50/70} ↑	ARCI _{50/70} ↓	ARTC _{50/70} ↑	AP _{50/70} ↑	C.V. ↓
mean		84.67/80.68	67.01/60.04	2.38/1.70	12.35/9.46	67.80/60.01	19.23

C. Performance analysis

C.1. Computation complexity

The configuration of the experiment platform has been described in Section B.1. Based on that, in terms of computations, the proposed entropy-cs only requires about 0.136G FLOPS with 0.66 ms latency to process a $256 \times 32 \times 32$ (C,H,W) feature map (more architecture details are shown in Section F.1).

Compared to DiscoNet[19], the proposed C-GRU costs about 4.10G FLOPS more with 12.7 ms latency.

C.2. Is Early Fusion always be better?

We discuss the performance of early fusion model. As we all know, Early fusion aggregates the raw measurements from all collaborators, promoting a holistic perspective. From the Table 1 in manuscript, the early fusion performs extremely good, even better than all the other baselines in some metrics. However, V2VNet[40] actually shows early fusion is far from optimal due to noises in real sensors. To address the above issue, since the dataset of V2VNet is not open source, we conduct experiments on OPV2V[46] with Gaussian noises. As shown in Table 8, we agree that the performance of Early Fusion may be degraded by noises to some extent.

Table 8. Comparisons on OPV2V dataset. [Best, Worst]

Method	ARSV _{50/70}	ARCV _{50/70}	AP _{50/70}
No Fusion	69.33/46.35	12.32/4.36	54.69/23.94
Early Fusion	64.62/46.95	45.00/24.39	55.88/25.89
UMC	76.56/47.68	47.82/25.06	61.90/24.50

C.3. Grains selection

Table 3 in manuscript compares the performance of different selections of grain level. We also include comparisons of single-grain, as shown in Table 9. Note that the heavy memory burden of all resolution baseline is not applicable on our RTX 3090.

Table 9. Comparisons of single-grain selection.

Multi-Grains Selection			AP _{50/70} ↑	C. V. ↓
$F_{i,1}^{e,t}$	$F_{i,2}^{e,t}$	$F_{i,3}^{e,t}$		
✓			56.73/49.22	7.88
	✓		58.96/52.86	8.12
		✓	57.43/51.66	8.36

D. More details about ablations analysis

Table 4 in manuscript shows that a tremendous drop when adding Entropy-CS in variant 3 and 4. From our perspectives, variant 3 and 4 are based on single-resolution, then variant 3 (w/ entropy-cs) costs about $\frac{1}{4}$ communication of variant 4. Based on [12], when the communication is too small, the collaborative detection performance will suffer, resulting in a tremendous drop in variant 3. However, variant 3 still achieves detection gain compared with No Fusion (improved by 9.40%/9.25% ↑ in AP_{50/70}, respectively).

Meanwhile, AP of variant 2 is worse than variant 4 with comparable ARSV and better ARCV, this is because The ARSV and ARCV only take recall into account, while the AP is weighted by both recall and precision. Therefore, in variants 2 and 4, a high AP does not necessarily mean high ARSV or ARCV values, more details about proposed metrics can be found in Section A.

E. Unified framework design

We summarize the main contributions of recent collaborative algorithms in Table 10, where ✓ indicates that a unique module is designed and – indicates that general operations are utilized. Our proposed UMC optimizes the communication, collaboration, and reconstruction process with multi-resolution technique.

F. Detailed architecture of the model

Note that we will release the source code.

F.1. Architecture of entropy-CS

Table 10. Contribution summary.

Method	Comm.	Collab.	Recons.
Who2com (ICRA 2020 [24])	✓	-	-
When2com (CVPR 2020 [23])	✓	-	-
V2VNet (ECCV 2020 [40])	-	✓	-
DiscoNet (NIPS 2021 [19])	-	✓	-
V2X-ViT (ECCV 2022 [45])	-	✓	-
Where2comm (NIPS 2022 [12])	✓	✓	-
UMC (ours)	✓	✓	✓

```

1  def acc_entropy_selection(self, tg_agent,
2     nb_agent, delta1, delta2, M=3, N=3):
3     self.stack = stack_channel(1, 9, kernel_size=3,
4     padding=1)
5     w = nb_agent.shape[-2]
6     h = nb_agent.shape[-1]
7     batch_nb = nb_agent.reshape(-1, 1, 1, 1)
8     stack = self.stack(nb_agent).permute(2,3,1,0).
9     contiguous().reshape(-1, 9, 1, 1)
10
11    p = F.sigmoid((stack - batch_nb)).mean(dim=1).
12    reshape(w, h)
13    entropy_tmp = p * torch.log(p)
14
15    with torch.no_grad():
16    top_delta = torch.sort(entropy_tmp.reshape(-1),
17    descending=True)
18    self_holder = top_delta[0][int(w*h*delta1)]
19
20    masker = torch.where(entropy_tmp >= self_holder)
21
22    stack_tg = self.stack(tg_agent).permute
23    (2,3,1,0).contiguous().reshape(-1, 9, 1, 1)
24    p_t = F.sigmoid((stack_tg - batch_nb)).mean(dim
25    =1).reshape(w, h)
26    entropy_t = p_t * torch.log(p_t)
27
28    tmp_masker = - torch.ones_like(entropy_t)
29    tmp_masker[masker] = entropy_t[masker]
30
31    with torch.no_grad():
32    top_delta2 = torch.sort(tmp_masker[tmp_masker
33    != -1].reshape(-1), descending=True)
34    thresholds = top_delta2[0][int(w*h*delta2)]
35
36    return torch.where(tmp_masker >= thresholds)
37
38    class stack_channel(nn.Conv2d):
39    def __init__(self, in_channels, out_channels,
40    kernel_size, stride=1, padding=0, bias=False,
41    interplate='none'):
42    super(stack_channel, self).__init__(in_channels
43    , out_channels, kernel_size=kernel_size,
44    stride=stride, padding=padding, bias=bias)
45
46    square_dis = np.zeros((out_channels,
47    kernel_size, kernel_size))
48
49    for i in range(out_channels):
50    square_dis[i, i//3, i%3] = 1
51
52    self.square_dis = nn.Parameter(torch.Tensor(
53    square_dis), requires_grad=False)
54
55    def forward(self, x):

```

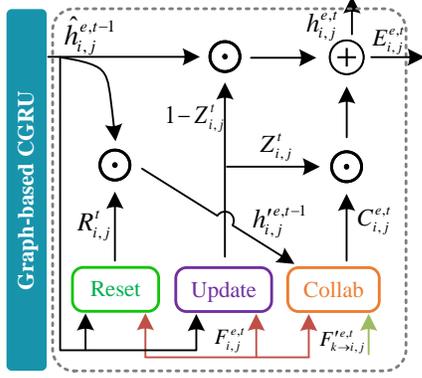


Figure 12. The architecture of G-CGRU.

```

44 kernel = self.square_dis.detach().unsqueeze(1)
45 stack = F.conv2d(x, kernel, stride=1, padding
46             =1, groups=1)
47 return stack

```

Listing 2. Entropy-CS code

Our contribution of entropy-based selection is both theoretical and practical. The intuition of entropy-cs is low computational complexity and high interpretability. The entropy-cs is no parameter and single-round communication to reduce heavy bandwidth burden brought by multi-resolution technique.

F.2. Architecture of G-CGRU

To facilitate understanding, we have simplified many formulas and steps in manuscript. Hence, we add more technique details about the section of Graph-based Collaborative GRU.

For the ego agent i of the j -th resolution intermediate feature maps, the inputs of G-CGRU are hidden states $\hat{h}_{i,j}^{e,t-1}$, the ego agent observation $F_{i,j}^{e,t}$, and the supporters' selected feature maps $\{F_{k \to i,j}^{e,t}\}_{k \neq i}$, then the updates for G-CGRU at t -th step can be formulated as:

$$\begin{aligned}
\hat{h}_{i,j}^{e,t-1} &= \Lambda(\mathbf{h}_{i,j}^{e,t-1}, \xi_i^{t-1 \rightarrow t}) \\
\mathbf{R}_{i,j}^t &= \text{Reset}(\mathbf{F}_{i,j}^{e,t}, \hat{h}_{i,j}^{e,t-1}) \\
\mathbf{Z}_{i,j}^t &= \text{Update}(\hat{h}_{i,j}^{e,t-1}, \mathbf{F}_{i,j}^{e,t}) \\
\mathbf{h}_{i,j}'^{e,t-1} &= \hat{h}_{i,j}^{e,t-1} \odot \mathbf{R}_{i,j}^t \\
\mathbf{C}_{i,j}^{e,t} &= \text{Collab}(\mathbf{h}_{i,j}'^{e,t-1}, \{\mathbf{F}_{k \to i,j}^{e,t}\}_{k \neq i}, \mathbf{F}_{i,j}^{e,t}) \\
\mathbf{E}_{i,j}^{e,t} &= \mathbf{Z}_{i,j}^t \odot \mathbf{C}_{i,j}^{e,t} + (1 - \mathbf{Z}_{i,j}^t) \odot \hat{h}_{i,j}^{e,t-1} \\
\mathbf{h}_{i,j}^{e,t} &= \mathbf{W}_{3 \times 3} * \mathbf{E}_{i,j}^{e,t}
\end{aligned} \tag{7}$$

where \odot , $*$ represent dot product and 3×3 convolution operation, respectively. $\mathbf{W}_{3 \times 3}$ indicates trainable parameters. The last time hidden feature $\mathbf{h}_{i,j}^{e,t-1}$ needs conduct

feature alignment operation from $t-1$ to t to get $\hat{h}_{i,j}^{e,t-1}$, which ensure the $\hat{h}_{i,j}^{e,t-1}$ and $\mathbf{F}_{i,j}^{e,t}$ are supported in the same coordinate system.

The *Reset* and *Update* gate modules share the same structure. Here we take *Reset* as an example:

$$\begin{aligned}
\mathbf{W}_{ir} &= \sigma(\mathbf{W}_{3 \times 3} * ([\hat{h}_{i,j}^{e,t-1}; \mathbf{F}_{i,j}^{e,t}])) \\
\mathbf{R}_{i,j}^t &= \sigma(\mathbf{W}_{ir} \odot \hat{h}_{i,j}^{e,t-1} + (1 - \mathbf{W}_{ir}) \odot \mathbf{F}_{i,j}^{e,t})
\end{aligned} \tag{8}$$

where $\sigma(\cdot)$, $[\cdot; \cdot]$ represent Sigmoid function and concatenation operation along channel dimensions. The gate $\mathbf{R}_{i,j}^t \in \mathbb{R}^{K,K,C}$ learns where the hidden features $\hat{h}_{i,j}^{e,t-1}$ are conducive to the present.

Based on the above *Reset* and *Update* modules, we thus derive the *Collab* module. To make better collaborative feature integration, we construct a collaboration graph $\mathcal{G}_c^t(\mathcal{V}, \mathcal{E})$ in *Collab* module, where node $\mathcal{V} = \{\mathcal{V}_i\}_{i=1, \dots, N}$ is the set of collaborative agents in environment and $\mathcal{E} = \{\mathbf{W}_{i \rightarrow j}\}_{i,j=1, \dots, N}$ is the set of trainable edge matrix weights between agents and models the collaboration strength between two agents. Let $\mathcal{C}_{\mathcal{G}_c^t}(\cdot)$ be the collaboration process defined in the *Collab* module's graph \mathcal{G}_c^t . The j -th resolution enhanced maps of ego i agent after collaboration are $\mathbf{E}_{i,j}^{e,t} \leftarrow \mathcal{C}_{\mathcal{G}_c^t}(\mathbf{h}_{i,j}^{e,t-1}, \mathbf{F}_{k \rightarrow i,j}^{e,t}, \mathbf{F}_{i,j}^{e,t})$. This process has two stages: message attention (S1) and message aggregation (S2).

$$\begin{aligned}
\mathbf{W}_{k \rightarrow i} &= \Pi([\mathbf{h}_{i,j}'^{e,t-1}; \mathbf{F}_{k \rightarrow i,j}^{e,t}; \mathbf{F}_{i,j}^{e,t}] \in \mathbb{R}^{K,K}) \\
\bar{\mathbf{W}}_{k \rightarrow i} &= \frac{e^{\mathbf{W}_{k \rightarrow i}}}{\sum_{m=1}^N e^{\mathbf{W}_{m \rightarrow i}}} \\
\mathbf{C}_{i,j}^{e,t} &= \sum_{m=1}^N \bar{\mathbf{W}}_{m \rightarrow i} \odot \mathbf{F}_{m \rightarrow i,j}^{e,t}
\end{aligned} \tag{9}$$

In the *message attention stage* (S1), each agent determines the matrix-valued edge weights, which reflect the strength from one agent to another at each individual cell. To determine the edge weights, we firstly get the conducive history information from hidden features by *Reset* gates through $\mathbf{h}_{i,j}'^{e,t-1} \leftarrow \hat{h}_{i,j}^{e,t-1} \odot \mathbf{R}_{i,j}^t$. Then, we utilize the edge encode Π to correlate the history information, the feature map from another agent and ego feature map; that is, the matrix-value edge weight from k -th agent to the i -th agent is $\mathbf{W}_{k \rightarrow i} = \Pi([\mathbf{h}_{i,j}'^{e,t-1}, \mathbf{F}_{k \rightarrow i,j}^{e,t}, \mathbf{F}_{i,j}^{e,t}] \in \mathbb{R}^{K,K}$, where Π concatenates three feature maps along the channel dimension and then utilizes four 1×1 convolutional layers to gradually reduce the number of channels from $3C$ to 1, more details are shown in Section F.6. Also, to normalize the edge weights across different agents, we implement a softmax operation on each cell of the feature map.

In the *message aggregation stage (S2)*, each agent aggregates the feature maps from collaborators based on the normalized matrix-valued edge weights, the updated feature map $C_{i,j}^{e,t}$ is utilized by $\sum_{k=1}^N \mathbf{W}_{k \rightarrow i} \circ \mathbf{F}'_{k \rightarrow i,j}{}^{e,t}$, where \circ represents the dot production broadcasting along the channel dimension.

Finally, the collaborative map is $E_{i,j}^{e,t} = \mathbf{Z}_{i,j}^t \odot C_{i,j}^{e,t} + (1 - \mathbf{Z}_{i,j}^t) \odot \mathbf{h}_{i,j}^{e,t-1}$. Note that the $\mathbf{Z}_{i,j}^t$ is generated by *Update gate* and \odot is the dot product. And, the hidden state is updated as $\mathbf{h}_{i,j}^{e,t} \leftarrow \mathbf{W}_{3 \times 3} * E_{i,j}^{e,t}$.

F.3. Architecture of shared encoder

We use the main architecture of MotionNet[42] as our shared encoder. The input BEV map's dimension is $(c, w, h) = (13, 256, 256)$. We describe the architecture of the encoder below:

```

1  nn.Sequential(
2  nn.Conv2d(13, 32, 3, stride=1, padding=1)
3  nn.BatchNorm2d(32)
4  nn.ReLU()
5  nn.Conv2d(32, 32, 3, stride=1, padding=1)
6  nn.BatchNorm2d(32)
7  nn.ReLU()
8  nn.Conv3D(64, 64, (1,1,1), stride=1)
9  nn.Conv3D(128, 128, (1,1,1), stride=1)
10 nn.Conv2d(32, 64, 3, stride=2, padding=1)
11 nn.BatchNorm(64)
12 nn.ReLU()
13 nn.Conv2d(64, 128, 3, stride=2, padding=1)
14 nn.BatchNorm(128)
15 nn.ReLU()
16 nn.Conv2d(128, 128, 3, stride=1, padding=1)
17 nn.BatchNorm(128)
18 nn.ReLU()
19 nn.Conv2d(128, 256, 3, stride=2, padding=1)
20 nn.BatchNorm(256)
21 nn.ReLU()
22 nn.Conv2d(256, 256, 3, stride=1, padding=1)
23 nn.BatchNorm(256)
24 nn.ReLU()
25 nn.Conv2d(256, 512, 3, stride=2, padding=1)
26 nn.BatchNorm(512)
27 nn.ReLU()
28 nn.Conv2d(512, 512, 3, stride=1, padding=1)
29 nn.BatchNorm(512)
30 nn.ReLU())

```

Listing 3. Shared encoder code

F.4. Architecture of SRAR-based shared decoder

The input of the SRAR-based shared decoder is the intermediate feature output by each layer of the encoder. Its architecture is shown below:

```

1  nn.Sequential(
2  nn.Conv2d(512 + 256, 256, 3, 1, 1)
3  nn.BatchNorm2d(256)
4  nn.ReLU()
5  nn.Conv2d(256, 256, 3, 1, 1)
6  nn.BatchNorm2d(256)
7  nn.ReLU()

```

```

8  nn.Conv2d(256 + 128, 128, 3, 1, 1)
9  nn.BatchNorm2d(128)
10 nn.ReLU()
11 nn.Conv2d(128, 128, 3, 1, 1)
12 nn.BatchNorm2d(128)
13 nn.ReLU()
14 nn.Conv2d(128 + 64, 64, 3, 1, 1)
15 nn.BatchNorm2d(64)
16 nn.ReLU()
17 nn.Conv2d(64, 64, 3, 1, 1)
18 nn.BatchNorm2d(64)
19 nn.ReLU()
20 nn.Conv2d(64 + 32, 32, 3, 1, 1)
21 nn.BatchNorm2d(32)
22 nn.ReLU()
23 nn.Conv2d(32, 32, 3, 1, 1)
24 nn.BatchNorm2d(32)
25 nn.ReLU()

```

F.5. Architecture of query generator

The entropy-CS compresses the intermediate feature to generate query matrix by query generator, which is for light communication. Its architecture is shown below:

```

1  nn.Sequential(
2  nn.Conv2d(256, 64, 1, 1, 0)
3  nn.BatchNorm2d(64)
4  nn.ReLU()
5  nn.Conv2d(64, 1, 1, 1, 0)
6  nn.ReLU())

```

Listing 4. Query generator code

F.6. Architecture of edge encoder II

```

1  class EdgeEncoder(nn.Module):
2  def __init__(self, channel):
3  super(EdgeEncoder, self).__init__()
4
5  self.conv1_1 = nn.Conv2d(channel, 128,
6  kernel_size=1, stride=1, padding=0)
7  self.bn1_1 = nn.BatchNorm2d(128)
8
9  self.conv1_2 = nn.Conv2d(128, 32, kernel_size
10 =1, stride=1, padding=0)
11 self.bn1_2 = nn.BatchNorm2d(32)
12
13 self.conv1_3 = nn.Conv2d(32, 8, kernel_size=1,
14 stride=1, padding=0)
15 self.bn1_3 = nn.BatchNorm2d(8)
16
17 self.conv1_4 = nn.Conv2d(8, 1, kernel_size=1,
18 stride=1, padding=0)
19
20 def forward(self, x):
21 x = x.view(-1, x.size(-3), x.size(-2), x.size
22 (-1))
23 x_1 = F.relu(self.bn1_1(self.conv1_1(x)))
24 x_1 = F.relu(self.bn1_2(self.conv1_2(x_1)))
25 x_1 = F.relu(self.bn1_3(self.conv1_3(x_1)))
26 x_1 = F.relu(self.conv1_4(x_1))
27
28 return x_1

```

Listing 5. Edge encoder code


```

6 self.guide_v1_bn = nn.BatchNorm2d(128)
7 self.guide = nn.Conv2d(256, 256, kernel_size=1,
8   stride=1, padding=0)
9 self.guide_bn = nn.BatchNorm2d(256)
10 self.conv5_1 = nn.Conv2d(512 + 256 + 256, 256,
11   kernel_size=3, stride=1, padding=1)
12 self.bn5_1 = nn.BatchNorm2d(256)
13 self.conv5_2 = nn.Conv2d(256, 256, kernel_size
14   =3, stride=1, padding=1)
15 self.bn5_2 = nn.BatchNorm2d(256)
16 self.conv6_1 = nn.Conv2d(256 + 128 + 128, 128,
17   kernel_size=3, stride=1, padding=1)
18 self.bn6_1 = nn.BatchNorm2d(128)
19 self.conv6_2 = nn.Conv2d(128, 128, kernel_size
20   =3, stride=1, padding=1)
21 self.bn6_2 = nn.BatchNorm2d(128)
22 self.conv7_1 = nn.Conv2d(128 + 64, 64,
23   kernel_size=3, stride=1, padding=1)
24 self.conv7_2 = nn.Conv2d(64, 64, kernel_size=3,
25   stride=1, padding=1)
26 self.bn7_1 = nn.BatchNorm2d(64)
27 self.bn7_2 = nn.BatchNorm2d(64)
28 self.bn8_1 = nn.BatchNorm2d(32)
29 self.bn8_2 = nn.BatchNorm2d(32)
30
31 self.norm1 = L2Norm(512)
32 self.norm2 = L2Norm(256)
33 self.norm3 = L2Norm(128)
34 self.norm4 = L2Norm(64)
35 self.norm5 = L2Norm(32)
36
37 def forward(self, x, x_1, x_2, x_3, x_4,
38   enhance_v1, enhance, batch, kd_flag = 0):
39
40   enhance_v1 = enhance_v1.view(batch, -1,
41     enhance_v1.size(1), enhance_v1.size(2),
42     enhance_v1.size(3))
43   enhance_v1 = enhance_v1.permute(0, 2, 1, 3, 4).
44     contiguous()
45   enhance_v1 = enhance_v1.permute(0, 2, 1, 3, 4).
46     contiguous()
47   enhance_v1 = enhance_v1.view(-1, enhance_v1.
48     size(2), enhance_v1.size(3), enhance_v1.size
49     (4)).contiguous()
50
51   guide_v1 = torch.max(F.relu(self.guide_v1_bn(
52     self.guide_v1(enhance_v1))), dim=1, keepdim=
53     True)[0]
54   guide = torch.max(F.relu(self.guide_bn(self.
55     guide(enhance))), dim=1, keepdim=True)[0]
56   x_3_guide = guide * x_3
57
58   x_5 = F.relu(self.bn5_1(self.conv5_1(torch.cat
59     ((self.norm1(F.interpolate(x_4, scale_factor
60     =(2, 2))), self.norm2(x_3_guide), self.norm2(
61     enhance))), dim=1))))
62   x_5 = F.relu(self.bn5_2(self.conv5_2(x_5)))
63
64   x_2 = x_2.view(batch, -1, x_2.size(1), x_2.size
65     (2), x_2.size(3))
66   x_2 = x_2.permute(0, 2, 1, 3, 4).contiguous()
67   x_2 = x_2.permute(0, 2, 1, 3, 4).contiguous()
68   x_2 = x_2.view(-1, x_2.size(2), x_2.size(3),
69     x_2.size(4)).contiguous()
70
71   x_2_guide = guide_v1 * x_2
72
73   x_6 = F.relu(self.bn6_1(self.conv6_1(torch.cat
74     ((self.norm2(F.interpolate(x_5, scale_factor
75     =(2, 2))), self.norm3(x_2_guide), self.norm3(
76     enhance_v1))), dim=1))))
77
78   x_6 = F.relu(self.bn6_2(self.conv6_2(x_6)))
79
80   x_1 = x_1.view(batch, -1, x_1.size(1), x_1.size
81     (2), x_1.size(3))
82   x_1 = x_1.permute(0, 2, 1, 3, 4).contiguous()
83   x_1 = x_1.permute(0, 2, 1, 3, 4).contiguous()
84   x_1 = x_1.view(-1, x_1.size(2), x_1.size(3),
85     x_1.size(4)).contiguous()
86
87   x_7 = F.relu(self.bn7_1(self.conv7_1(torch.cat
88     ((self.norm3(F.interpolate(x_6, scale_factor
89     =(2, 2))), self.norm4(x_1), dim=1))))
90   x_7 = F.relu(self.bn7_2(self.conv7_2(x_7)))
91
92   x = x.view(batch, -1, x.size(1), x.size(2), x.
93     size(3))
94   x = x.permute(0, 2, 1, 3, 4).contiguous()
95   x = x.permute(0, 2, 1, 3, 4).contiguous()
96   x = x.view(-1, x.size(2), x.size(3), x.size(4))
97     .contiguous()
98
99   x_8 = F.relu(self.bn8_1(self.conv8_1(torch.cat
100     ((self.norm4(F.interpolate(x_7, scale_factor
101     =(2, 2))), self.norm5(x), dim=1))))
102   res_x = F.relu(self.bn8_2(self.conv8_2(x_8)))
103
104   return res_x
105
106 class L2Norm(nn.Module):
107   def __init__(self, n_channels, scale=10.0):
108     super(L2Norm, self).__init__()
109     self.n_channels = n_channels
110     self.scale = scale
111     self.eps = 1e-10
112     self.weight = nn.Parameter(torch.Tensor(self.
113       n_channels))
114     self.weight.data *= 0.0
115     self.weight.data += self.scale
116
117   def forward(self, x):
118     norm = x.pow(2).sum(dim=1, keepdim=True).sqrt()
119     + self.eps
120     x = x / norm * self.weight.view(1, -1, 1, 1)
121
122   return x

```

Listing 6. MGFE code

G. Detailed information of Interpolation

We utilize the main architecture of ADP-C[25] as our interpolate function in entropy-CS module. We assume the input feature as $\mathbf{f}_{in} \in \mathbb{R}^{K,K,C}$ and suppose the pixels of the \mathbf{f}_{in} are indexed by \mathbf{p} . Then, we form a mask $\mathbf{M} \in \mathbb{R}^{K,K}$:

$$\mathbf{M}(\mathbf{p}) = \begin{cases} 0, & \text{if } \mathbf{f}_{in}(\mathbf{p}) = \mathbf{0} \\ 1, & \text{otherwise} \end{cases} \quad (10)$$

Assuming C is a convolution layer with input \mathbf{f}_{in} , the by applying the mask, the output \mathbf{f}_{out} at position \mathbf{p} becomes:

$$\mathbf{f}_{out}(\mathbf{p}) = \begin{cases} C(\mathbf{f}_{in})(\mathbf{p}), & \text{if } \mathbf{M}(\mathbf{p}) = 1, \\ \mathbf{0}, & \text{if } \mathbf{M}(\mathbf{p}) = 0. \end{cases} \quad (11)$$

Denoting the interpolation operation as I , the final output feature \mathbf{f}_{out}^* is:

$$\mathbf{f}_{out}^*(\mathbf{p}) = \begin{cases} \mathbf{f}_{out}(\mathbf{p}), & \text{if } \mathbf{M}(\mathbf{p}) = 1, \\ I(\mathbf{f}_{out})(\mathbf{p}), & \text{if } \mathbf{M}(\mathbf{p}) = 0. \end{cases} \quad (12)$$

The value of $I(\mathbf{f}_{out})(\mathbf{p})$ is weighted average of all the neighboring pixels centered at \mathbf{p} within a radius r :

$$I(\mathbf{f}_{out})(\mathbf{p}) = \frac{\sum_{\mathbf{s} \in \Psi(\mathbf{p})} \mathbf{W}_{(\mathbf{p},\mathbf{s})} \mathbf{f}_{out}(\mathbf{s})}{\sum_{\mathbf{s} \in \Psi(\mathbf{p})} \mathbf{W}_{(\mathbf{p},\mathbf{s})}} \quad (13)$$

where \mathbf{s} indicates \mathbf{p} 's neighboring pixels and $\Psi(\mathbf{p}) = \{\mathbf{s} \mid \|\mathbf{s} - \mathbf{p}\|_{\infty} \leq r, \mathbf{s} \neq \mathbf{p}\}$, the neighborhood of \mathbf{p} . In UMC, we set radius $r = 7$. $\mathbf{W}_{(\mathbf{p},\mathbf{s})}$ is the weight assigned to point \mathbf{s} for interpolating at \mathbf{p} , for which we utilize the RBF kernel, a distance-based exponential decaying weighting scheme:

$$\mathbf{W}_{(\mathbf{p},\mathbf{s})} = \exp(-\lambda^2 \|\mathbf{p} - \mathbf{s}\|_2^2) \quad (14)$$

with λ being a trainable parameter. This indicates that the closer \mathbf{s} is to \mathbf{p} , the larger its assigned weight will be. Note that masked-out features $\mathbf{M}(\mathbf{p}) = 0$ still participate in the interpolation process as inputs with values of $\mathbf{0}$.

H. Detailed Qualitative results

We visualize the detection results between different collaborative approaches on V2X-sim[18] and OPV2V[46] datasets, as shown in Figure 13 and 14.

I. Loss curve

We visualize the loss curve of UMC, Early Fusion, When2com, Where2comm, V2VNet and DiscoNet on V2X-Sim and OPV2V on Figure 15 and 16, respectively.

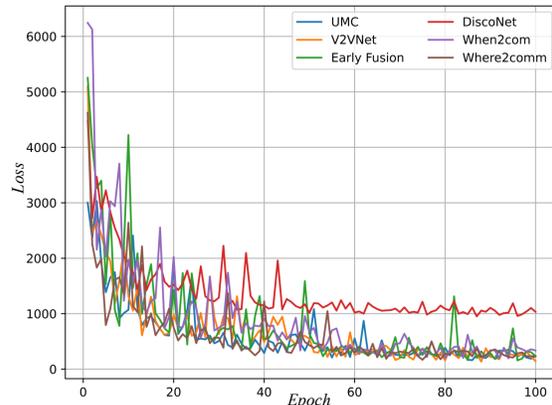


Figure 15. Epoch vs. loss on V2X-Sim dataset.

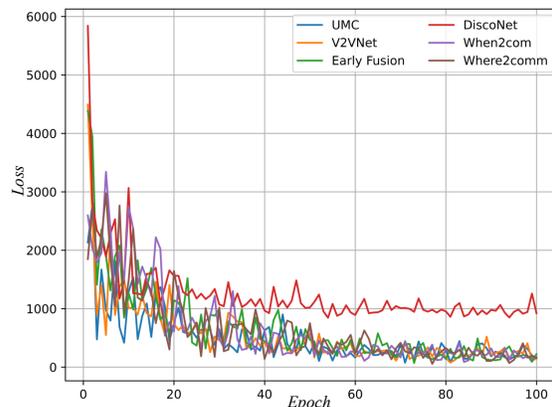


Figure 16. Epoch vs. loss on OPV2V dataset.