

Self-Attention Factor Graph Neural Network for Multi-Agent Collaborative Target Tracking

Cheng Xu, *Member, IEEE*, Ran Su, Ran Wang, *Graduate Student Member, IEEE*, and Shihong Duan, *Member, IEEE*

Abstract—Collaborative target tracking is an essential task in positioning systems, particularly in environments characterized by high dynamics, multi-source heterogeneous data, and interactive multi-agent scenarios. The challenge in such networks lies in the direct utilization of multi-source heterogeneous data as feature input for models. Additionally, the presence of high-dynamic time series data complicates the extraction of dependencies by the models. To address these issues, we introduce a novel approach that integrates a factor graph-based data fusion method with a graph neural network. This combination is designed to uncover potential dependencies between time series data and positional information within dynamic networks. Furthermore, we employ a self-attention mechanism, enabling distance-agnostic autonomous selection of complex network features. This innovation allows the model to achieve enhanced accuracy performance while simultaneously reducing computational costs. We validated our approach through simulation experiments. The results demonstrated the method's effectiveness in fusing and selecting multi-source heterogeneous information within collaborative networks. It also excelled in identifying potential relationships between feature information and positional data, showcasing the robustness and applicability of our proposed solution in challenging collaborative target tracking environments.

Index Terms—collaborative tracking, multi-agent network, factor graph, graph neural network, factor graph neural network.

I. INTRODUCTION

WITH the advent of the Internet of Things (IoT), accurate localization has become a cornerstone in the labyrinth of smart devices that pervade our daily lives. While the Global Positioning System (GPS) remains the bedrock of outdoor localization due to its high precision and cost-efficiency, its effectiveness is markedly reduced indoors, where satellite signals are scarce and prone to obstruction. In the quest for reliability within indoor and complex environments, alternative wireless positioning technologies such as Time-of-Arrival (TOA) and Time-Difference-of-Arrival (TDOA) have risen to prominence [1], [2]. These methods, however, grapple with

occlusions and non-line-of-sight conditions that are endemic to cluttered, dynamic collaborative networks [3]. Such conditions frequently degrade the accuracy of positioning, presenting a significant challenge in densely packed IoT ecosystems.

Compounding these challenges are the limitations of traditional inertial systems. While they offer notable instantaneous accuracy at a lower cost, their reliance on integration means that any drift in accelerometer or gyroscope readings can lead to increasing errors over time [4]. The inherent multi-source heterogeneity of sensor data in dynamic IoT networks further complicates the landscape, posing a formidable barrier for traditional positioning methodologies to harness this wealth of information effectively.

In complex and highly dynamic collaborative networks, achieving effective collaborative distribution, data screening, and fusion poses significant challenges in network positioning. The method of multi-source heterogeneous information fusion effectively harnesses data from various sensors, circumventing environmental or target perception limitations, thereby enhancing the system's external perception capabilities [5]. Wymeersch et al. [6] implemented a Bayesian approach based on message propagation for processing multi-source heterogeneous data to determine agent locations. However, the application of this method to highly dynamic networks, especially those involving time series data, is limited due to its high computational demand and the prerequisite of prior knowledge, such as offline labeling and positioning networks [7]. To transcend these constraints in current location technologies, there is a pressing need for a method capable of fusing multi-source data and incorporating time series information.

In recent years, graph-structure-based algorithms have gained widespread use in various graph-related learning tasks, including node classification, link prediction, and graph classification [8]–[11]. Concurrently, the use of graph structures has been extended to positioning [7]. Graph structures, compared to traditional methods, can uncover additional node interrelations, thereby revealing more potential interaction information. In the domains of artificial intelligence and neural networks, stochastic models are typically represented as Bayesian networks or Markov random processes [12]. However, in real-world scenarios, time series data often constitute incomplete observations of intricate underlying dynamic processes with high-dimensional states that are not directly observable. For instance, human motion capture data provide specific marker positions, reflecting the complex kinematic and dynamic constraints of numerous joint angles [13]. Traditional graph structures, however, often fall short in depicting these dependencies.

Manuscript received XX 2023; revised XX 2024; accepted XX 2024. Date of publication XX 2024; date of current version XX 2024. This work is supported in part by National Natural Science Foundation of China (NSFC) under Grant 62101029, and in part by the China Scholarship Council Award under Grant 202006465043. (*Corresponding author: Shihong Duan*)

The authors are with School of Computer and Communication Engineering, University of Science and Technology Beijing, Beijing, China. They are also with Shunde Innovation School, University of Science and Technology Beijing, Foshan, China. (email: xucheng@ustb.edu.cn; 2205372000@qq.com; wangran423@foxmail.com; duansh@ustb.edu.cn).

Copyright (c) 20xx IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

Multi-agent localization is a multifaceted process that encompasses the interaction between agents and the potential unobstructed variable information of these agents [14]. Even with the integration of multi-source data, extracting useful information from large datasets remains a significant challenge. Mirowski et al. [13] introduced a factor graph approach tailored for time series data to elucidate the relationship between time series and hidden states. Unlike traditional graph structures, factor graphs are extensively utilized to model independent random variables. For instance, approximate inference algorithms are applied in probabilistic graph models to reason about specified probability graphs [15], enabling a more accurate representation of the relationships between state variables and observable variables. In collaborative positioning scenarios that include timing information, factor graphs can ascertain higher-order potential dependencies, thereby facilitating the analysis of correlations between time series and data features [16]. However, in multi-agent cooperative networks, defining an effective prior model is challenging, often resulting in suboptimal approximate inference outcomes. Thus, efficiently processing data features in highly dynamic complex networks to enable automatic learning of inferential relationships and feature screening poses a formidable challenge.

Addressing this issue, Gori et al. [17] proposed the graph neural network (GNN) method, which learns potential variables and reasoning processes from data. This method only requires the provision of a graph structure with dependency relationships to better explore information within real data. Moreover, Zhang et al. [16] developed the factor graph neural network method, extending factor graphs to graph neural networks. This approach effectively tackles point cloud data classification but still struggles with the screening of complex input data. Velickovic et al. [10] introduced a graph attention method that filters data by assigning weights to different graph nodes. However, this method requires recalculating graph attention when the graph structure changes, making it less suitable for the real-time demands of dynamic networks.

Furthermore, self-attention mechanisms have become crucial in various tasks, enabling the modeling of dependencies regardless of the distance between input and output sequences [18]–[20]. In terms of computational complexity, self-attention offers parallel computation advantages compared to networks like RNNs. This mechanism processes the input sequence into queries, keys, and values, subsequently deriving a weight coefficient by calculating the similarity between the query and the key. The final output is a weighted summation of the values and these weight coefficients. The parallel nature of this computation primarily lies in the similarity calculation between queries and keys, achieved through matrix multiplication. Since matrix multiplication is highly parallelizable, it significantly accelerates model computation. When the sequence length n is less than the representation dimension d , the self-attention layer outperforms recurrent layers, making it more efficient for short-sequence timing information [20]. In practical positioning scenarios, we often use timing information close to the current moment for computation. Therefore, the self-attention mechanism can address data screening challenges while reducing computational costs in dynamic networks.

In this paper, we propose a collaborative localization model based on factor graph self-attention. This model leverages the synergistic integration of self-attention mechanisms and factor graphs within a graph neural network framework, enabling the effective screening of features from multi-agent nodes. This approach addresses the challenges of multi-source heterogeneous data fusion and screening, and realizes multi-agent localization through the use of factor graph neural networks to learn the reasoning process of data features. The primary contributions of our research are as follows:

- 1) *General Factor Graph-Based Framework for Collaborative Target Tracking*: We introduce a comprehensive framework that incorporates factor graphs within the realm of graph neural networks. This framework is specifically designed to navigate the unique challenges found in collaborative target tracking, such as dynamic environments and intricate agent interactions. It marks a significant advancement over conventional cooperative network approaches by enabling adaptive learning and meeting the real-time demands of dynamic scenarios.
- 2) *Advanced Data Fusion Method Utilizing Factor Graphs*: Our methodology presents an innovative data fusion approach, transforming heterogeneous data from a variety of sources into a structured graph format. This facilitates efficient processing of dynamic and complex datasets and enables the automatic discovery of complex inferential relationships, thus surpassing the constraints typically associated with traditional graph models.
- 3) *Integration of Self-Attention Mechanism in Factor Graph Neural Networks*: We have successfully integrated a self-attention mechanism into the factor graph neural network architecture. This enhancement bolsters the model's proficiency in processing timing information from agents and efficiently filtering interaction features. Consequently, our model adeptly extracts pertinent data features within highly dynamic and complex network environments, assisting collaborative agents in actively selecting perceptual data to optimize their pose estimation. This functionality is particularly beneficial in scenarios with uncertain network conditions.

The remainder of this paper is organized as follows: Section II introduces the related work and research, elaborating in detail on the merits of this study. Section III describes the localization method of factor graph neural networks based on self-attention, followed by experimental results and discussions in Section IV. Section V concludes the paper with a comprehensive summary.

II. RELATED WORK

A. Cooperative localization

Collaborative localization involves both the self-observation of agents and the measurement of interactive information between them. Current methods for localization primarily fall into two categories: optimization-based and learning-based approaches. Optimization-based methods include maximum likelihood estimation [21], [22], the least squares method [6], multidimensional scaling [23]–[25], and Bayesian message

propagation [6], [26], [27], among others. A comparative analysis of the methods presented in this paper and the pros and cons of existing methods are delineated in Table 1.

The maximum likelihood estimation method treats noise as a definitive probability distribution [28], [29], but this can lead to significant performance degradation in the event of incorrect matching [7]. In the realm of machine learning, Xie Hong et al. [30] developed an indoor positioning method that combines random forest and deep learning. This method employs deep learning to train the channel propagation model offline and determine the orientation of agents during the online stage. Yan et al. [7] introduced a location method based on graph neural networks (GNN), which addresses the computational performance challenges in large-scale networks and provides enhanced accuracy and stability for static localization of large-scale agents. However, in highly dynamic cooperative networks where the agent target is typically dynamic and includes time series of historical moments, the GNN method struggles to effectively utilize time series information.

Addressing this gap, Liu, Zhijun et al. [31] proposed a dynamic representation learning framework for network embedding on large-scale attributed networks. This approach models time-varying network features for dynamic attribute network embedding. In the offline learning stage, node embedding is generated by discovering potential node attributes and network structure to guide the learning subspace. In the online learning stage, dynamic network features are processed through timely and incremental updating of node embedding. In terms of heterogeneous data, X. Wang et al. [32] introduced the HAN method, incorporating the attention mechanism into the heterogeneous graph neural network. This method maps node features of different dimensions to a uniform dimension, effectively addressing the challenge of heterogeneous data processing and providing a clear explanatory framework. Building upon this, our paper adds a self-attention method to the factor graph neural network framework, thereby effectively resolving the heterogeneous data processing challenge in the cooperative localization environment.

Addressing the complexities of high-dynamic environments, Partwardhan et al. [33] introduced a planning method based on confidence propagation, delving into robot planning challenges in high-speed, congested traffic settings. This method particularly addresses the issue of heterogeneous data in dynamic environments. By categorizing complex data from dynamic settings into distinct factor nodes according to their sources, the approach facilitates effective path planning. To capture long-term dynamic dependencies effectively, it's imperative for a model to maintain an internal state that adheres to dynamic constraints. To address this requirement, the concept of a dynamic factor graph has been proposed. In this context, factor graphs are applied to time series data to map out the dependency relationships between state variables and observable variables, thereby unveiling higher-order dependency relationships. Building on this foundation, Murai et al. [34] proposed a dynamic factor graph localization method that utilizes Gaussian confidence propagation. They developed a methodology that employs probabilistic factor graphs for perception and state estimation. By adopting Gaussian con-

TABLE I: The strengths and limitations of various kinds of work in collaborative target tracking applications.

Methods	Strengths	Limitations
TOA	High precision	High cost, easy to block interference
MLE	Easy to implement	Noise probability distribution matching results in performance degradation
GNN	Large scale, high stability	The dynamic environment is not satisfied
BP	Obtain high level dependencies for dynamic environments	Require prior knowledge
Ours	Obtain high order dependencies in dynamic environments with high accuracy and low cost	Certain agent scale

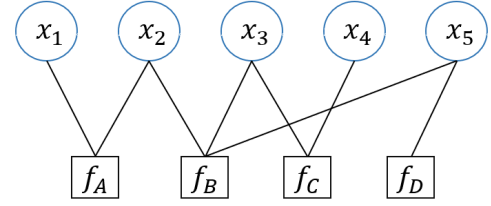


Fig. 1: Factor Graph

fidence propagation as the inference algorithm, this method enables agents to accurately estimate and adapt to dynamically collaborative scenes through self-organizing communication.

B. Factor Graph

Factor graphs are a class of graphical models used to infer tasks in probabilistic models by using graphs to represent dependencies between variables. The factor graph $G = (V, F, A)$ of bipartite graph is defined, which contains a group of variable nodes V , a group of factor nodes F , and a group of undirected edge A representing the connection relationship between variable nodes and factor nodes. Variable nodes represent hidden random variables and are usually represented by circles. Factor nodes represent probability-based dependencies between variable nodes, usually in square form. In the graph, every node $i \in V$ is associated with a random variable x_i , and every node $j \in F$ is associated with a function f_j . There is an edge connecting nodes f_i to nodes x_i if and only if the factor node f_i depends on the variable node x_i . For example, a function $f(x_1, x_2, x_3, x_4, x_5)$ with the sum of variable nodes x_1, x_2, x_3, x_4 and x_5 , factors the function into a product form

$$f(x_1, x_2, x_3, x_4, x_5) = f_A(x_1, x_2) f_B(x_2, x_3, x_4) f_C(x_3, x_4) f_D(x_5) \quad (1)$$

Wherein, it includes factor node $F = \{f_A, f_B, f_C, f_D\}$ and variable node $V = \{x_1, x_2, x_3, x_4, x_5\}$.

Over the past few decades, factor graphs have gained widespread use in coding and stochastic modeling for defining probability graph models and modeling dependencies between random variables. To address coding challenges, [35] introduced the sum-product algorithm based on belief propagation. This algorithm effectively solved the decoding problem of

LDPC codes, markedly enhancing their error correction capabilities. The sum-product algorithm [36], a message-passing algorithm, computes marginal distributions through a series of locally computed messages. It is adept at accurately or approximately calculating global functions [37]. Tanner [38] employed graphs to describe and generalize this encoding method, introducing the minimum sum algorithm. In recent years, factor graphs have emerged as an effective method for multi-source information fusion in collaborative localization. For instance, Von Stumberg et al. [39] proposed a factor graph-based approach that addresses computational redundancy and asynchronous delay in multi-source information processing, offering the advantage of plug-and-play functionality.

However, a significant challenge in practical scenarios is that we often obtain only an approximate value of the true distribution of the probability graph model, resulting in suboptimal outcomes. Additionally, message propagation technology, integral to factor graph methodologies, encounters limitations within cyclic graph structures. It is generally limited to calculating approximate values of the ideal posterior distribution, which can impede algorithm convergence [40]. Consequently, there has been growing interest in combining message propagation technology with graph networks in the coding field and classification tasks [16], [33], [40]. Satorras et al. [40] proposed a method that merges traditional message propagation with graph neural networks, applying it to LDPC coding to address complex noise issues in Gaussian channels. Additionally, Zhang et al. [16] integrated the factor graph model with graph networks for classifying point cloud data. By utilizing graph network algorithms, models can learn potential variables and inference processes from data, requiring only the graph structure to provide variable dependencies. While factor graphs excel at capturing higher-order potential dependencies in scenarios involving time series information, their direct application to target tracking is hindered by the common necessity to approximate the true distribution of the probability graph model in real-world scenarios, leading to less than optimal results. Moreover, the message propagation technique faces considerable limitations in cyclic graph structures, typically computing only approximate values of the ideal posterior distribution, which may prevent algorithm convergence and thus diminish the effectiveness of factor graph approaches in meeting the specific requirements of target tracking. Therefore, the integration of factor graph models with graph networks can enable models to actively learn the reasoning process of agents, discover potential motion features from time series and agent characteristics, and consequently achieve more accurate target localization.

C. Attention Mechanism

In the domain of collaborative positioning, the complexity of a vast array of input data features and the accumulation of time series information pose challenges in computational performance and associated costs for models. To address these challenges effectively, the attention mechanism has emerged as a standard in many sequential tasks. Its key advantage is the ability to process inputs of any size and focus on the most rele-

vant features within those inputs [10]. The self-attention mechanism connects information from different locations within the same sequence, facilitating the computation of sequence representations. This approach has seen successful applications in various tasks, including reading comprehension, summary generation, and text implication, by enabling the selection of effective features from massive datasets, thus enhancing model performance while reducing computational costs [41]–[44].

Graph attention mechanisms have gained widespread use in graph structures. Velickovic et al. [10] introduced an attention mechanism that filters data by assigning weights. However, the incorporation of double computation and time series in dynamic graph structures tends to increase computation time. Methods such as those proposed by Kosaraju et al. [45], Alahi et al. [46], and Gupta et al. [47] initially employ the Long Short-Term Memory model [48] to capture each agent's track features over time. These features are then fed into an interactive model graph neural network [8] to identify interactions between agents. However, such methods can lead to computational resource wastage, and issues like gradient vanishing and exploding can adversely impact the model's performance and accuracy.

To overcome these challenges, Schwartz et al. [49] proposed applying the graph attention mechanism to factor graphs for classifying image-based dialogue. This innovative approach has paved new paths for utilizing attention mechanisms within factor graphs. Inspired by these methods, this paper applies the self-attention mechanism to the factor graph network. This application enables the screening and extraction of features from accumulated time series data and the selection of effective features from numerous inputs, thereby enhancing the performance of the collaborative localization model.

III. SELF-ATTENTION FACTOR GRAPH NEURAL NETWORK

In this section, we define the symbolic definition of data items in the collaboration localization scenario, and then the data fusion of multi-source heterogeneous data in the highly dynamic network is introduced based on the factor graph. Factor graphs help us capture dependencies between information. We further extend the factor graph into factor graph neural network, and realize the autonomous feature screening and learning inference process of agents based on self-attention, so as to achieve cooperative localization.

A. Problem formulation

This paper assumes that the collaboration scenario is a two-dimensional wireless network space. Assume the current time is T ($T > 0$), $S_a = \{1, 2, \dots, N\}$ represents the directory of the agent collection, where N is the number of agents. $X^t = \{x_1^t, x_2^t, \dots, x_N^t\}$ represents the joint state of all N agents in t ($t < T$), including agent position p , velocity v , orientation Angle a and inter-agent distance data d ; $x_i^t = \{c_1, c_2, \dots, c_M\}$ represents agent measurement data at time T , including agent velocity v , orientation Angle a and distance between agents d at the current time. Let $Y^T = \{y_1^T, y_2^T, \dots, y_N^T\}$ represents the predicted state of all agents at time T . The objective of this paper is to locate the position information of all

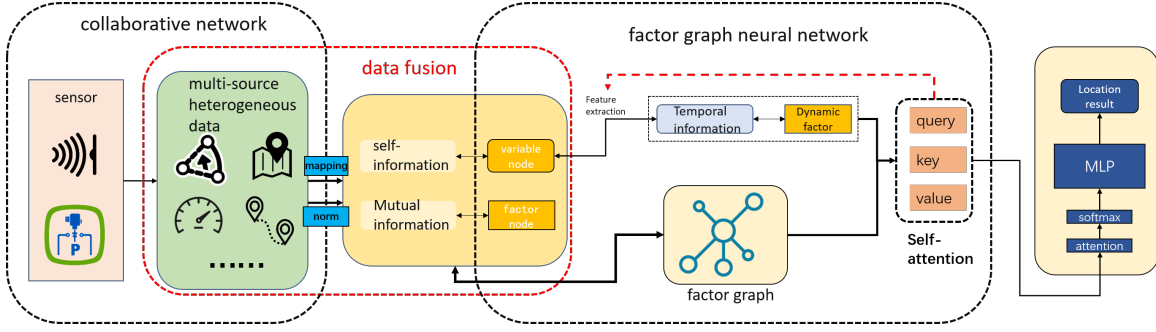


Fig. 2: Self-attention Factor Graph Neural Network.

agents at the current moment based on the measurement data of multiple agents at the current moment X^T and the observation state of historical moment in the collaborative network $X_{his} = \{X^t\}_{t=0}^{T-1}$, expressed as $\rho_\sigma(Y^{T+1}|X_{his}, x^T)$.

B. Multi-source data fusion

Aiming at the problem that multi-source heterogeneous data of agents cannot be directly used as model input in collaborative scenarios, this paper proposes a solution. First, we define the interaction characteristics of agents as self-information factor nodes, the characteristics of agents themselves as variable nodes, and the hidden feature information obtained based on self-attention time series data as dynamic factor nodes. Then, we set up the form of factor graph, fused multi-source heterogeneous data, and carried out feature mapping. by this means, we transform the raw data into agent characteristics acceptable to the model. In this way, we can effectively process multi-source heterogeneous data in collaborative scenarios and improve the accuracy and reliability of the model.

For the measurement data $x_i^T = \{c_1, c_2, \dots, c_M\}$ obtained by the sensor at time T of the multi-agent, according to the data dependency relationship, it is defined as the measurement data between agents obtained by communication sensors, namely mutual information X_{mut}^T , and the measurement data of the agent itself obtained by the sensor inside the agent, namely self-information X_{self}^T . Where, the self-information $X_{self}^T = \{X_i\}_{i=0}^N$ only represents the data of agent i , including the acceleration information and orientation Angle of agent i at time T , which is defined as a variable node in this paper. Mutual information $X_{mut}^T = (X_{ij})_{i=0, j=0}^{i,j=N}$ contains the ranging information between agent i and agent j at time T , which is dependent on the data of agent ij at a single time, which is defined as an internal factor node in this paper. The original data obtained by these sensors cannot be calculated directly as the input of the model due to the heterogeneity of the data resulting in different feature dimensions. Therefore, for different kinds of feature data c_i in x_1^t , this model uses multiple feature mapping layers φ_i and standardized functions θ_i to transform data features of different dimensions into the same d -dimension feature vector C_i , and then uniformly process them through standardized functions. Finally, all feature vectors are mapped to a d -dimensional feature space for calculation, as shown in Figure 3.

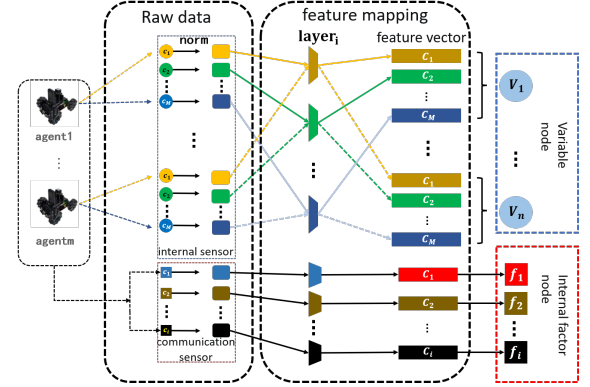


Fig. 3: For multi-source heterogeneous raw data, different feature mapping layers are used to transform it into standard input features, and the data sources are classified into different nodes of factor graph.

Sensor data c_i in collaborative network contains speed information, angle information and distance information respectively. According to the data dimensions of different structures, this paper defines the feature mapping function:

$$\varphi_i \leftarrow \phi_i(w_i^{D_i \times d}, d) \quad (2)$$

where $w_i^{D_i \times d}$ is the network parameter. Feature mapping function can transform multi-source heterogeneous data into unified dimension input features, which is convenient for model calculation. However, due to the different input data units and orders of magnitude of different structures, direct input will lead to poor performance of the model. In order to deal with this problem, this paper obtained k samples in simulation scenarios for data in collaborative scenarios, and defined standardized functions using Gaussian distribution:

$$\theta_i \leftarrow \frac{Q - \sum_{i=1}^k Q_i}{\sqrt{k} \times \sqrt{\sum_{i=1}^k (Q_i - \frac{\sum_{i=1}^k Q_i}{k})^2}} \quad (3)$$

where Q_i is sample data. According to the defined feature mapping function and standardized function, the multi-source heterogeneous data fusion function of this model as followed:

$$C_i = \theta_i(\varphi_i(c_i, w_i^{D_i \times d}, d)) \quad (4)$$

According to different data dependencies, this model takes the interaction features between agents extracted from data as the internal factor node F_i of the factor graph, and takes the features of agents themselves as the variable node V of the factor graph. In addition, the model represents the communication relationship between agent i and agent j as an adjacency matrix A_{ij} .

$$A_{ij} = \begin{cases} 0, & \text{no communication between agents } ij \\ 1, & \text{communication between agents } ij \end{cases} \quad (5)$$

Then the factor graph at time t is expressed as $G = (V, F_i, A)$. Such representation makes it easier for the model to deal with the communication relationship between agents, so that it can calculate and predict more accurately.

In the cooperative network of agents, the position of agents at the next moment may be affected by the historical motion state of agents. Therefore, in order to obtain the motion hiding state of the agent from the historical characteristics of the agent, we introduce dynamic factor nodes into the factor graph. This node connects the factor graph variable node of the current moment and the agent variable node of the historical moment to obtain the characteristics of the historical state of the agent. Specifically, for the agent characteristics $c_{his} = \{c^{T-k+1}, c^{T-k+2}, \dots, c^T\}$ before the observed time T , a variable time length $p(p > 0)$ is selected to obtain the agent historical state characteristics $c_p = \{c^{T-p+1}, c^{T-p+2}, \dots, c^T\}$, where $c^T = \{c_1^t, c_2^t, \dots, c_N^t\}$ represents the historical data characteristics of all N agents at the time t . We can get the historical feature matrix mat_{his} by the column matrix of all the features of historical moments c^T . In order to learn the relationship between historical states and hidden motion states, this paper adopts self-attention network to obtain hidden state S :

$$S = \text{self_attention}(\mu_i) = \text{softmax}\left(\frac{\mu_i^T \mu_i}{\sqrt{d}}\right) \mu_i \quad (6)$$

$$\mu_i^{(i \in \{1,2,3\})} = W_i(\omega_i, mat_{his}) \quad (7)$$

where W_i is trainable network, ω_i is the parameter of attention network, and S is the hidden state of timing information of variable nodes.

In Part B, we elaborate on the construction and utility of the factor graph within our algorithm, highlighting its contributions to enhancing data processing efficacy and model accuracy. The factor graph encapsulates the dynamic factor nodes, which correlate with the variable nodes representing the model's hidden states S , across adjacent time points, thus forming a continuous time factor graph. Through this construction, our model transforms cooperative network data into input features that are conducive to the model's requirements, capturing essential feature dependencies. This factor graph creation process not only refines the data representation but also optimizes the model's performance.

C. The establishment of factor graph neural network

Through the above multi-source heterogeneous data fusion process, we can establish a factor graph with continuous

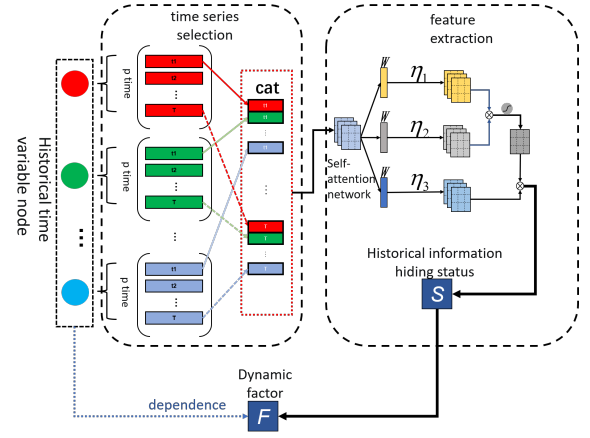


Fig. 4: For the obtained variable nodes, the hidden features of the history state are extracted from the attention network.

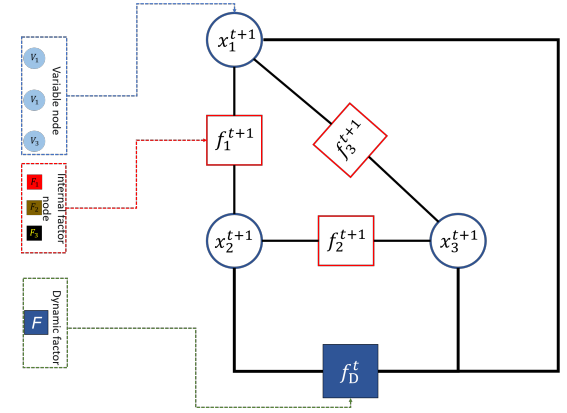


Fig. 5: Established factor diagram.

moments and obtain the connection dependence relationship between nodes. However, it is often difficult to obtain a priori model of feature inference in highly dynamic collaborative networks. Therefore, we introduce self-attention neural network to learn the feature inference relationship, and enable agents to actively select and screen the most beneficial feature information.

Figure 5 introduces the process of establishing factor graph in this paper by taking three agents as an example. For variable node V and internal factor F_i node obtained through information fusion at time t in the previous section, x_i^t is used in this paper to represent characteristic information of variable node V of agent i at time t . f_{ij}^t represents the internal factor node F between agent i and agent j at time t , represents that the variables x_i^t and x_j^t are dependent. Then there is an edge connecting the variable node V with the internal factor node F_i . For the dynamic factor node F_D obtained from the timing information of variable node, In this paper, f_D^t represents the dynamic factor of the agent obtained according to the historical characteristic state at time t , which is dependent on variable nodes of the factor graph at adjacent moments. Then, there is an edge connecting all variable nodes V and dynamic factor nodes F_D .

For factor graph $G = (V, F_i, F_D, A)$. In order to obtain

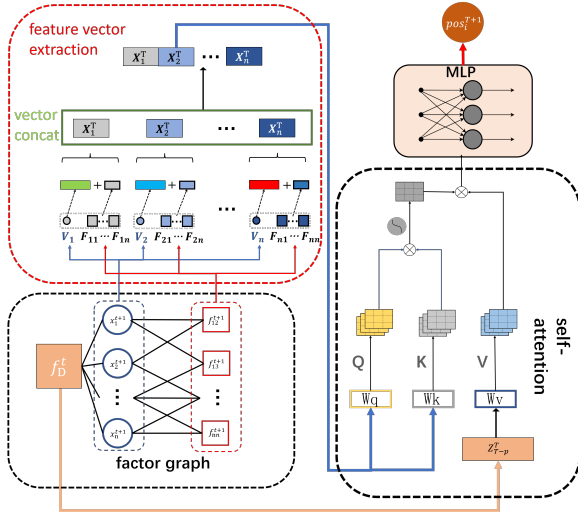


Fig. 6: By establishing factor graph neural network and based on self-attention, cooperative localization is realized.

Algorithm 1 Training

Input: $D_i (i \in M), c_i \in x_i^t, Q_i$, Time information X_{his} and label Y^T

for $i = M, \dots, 1$ **do**

 Create φ_i by (2).

 Create θ_i by (3).

 GET $(V_c, F_k) = C_i$ by (4).

end for

Obtain the adjacency matrix by

$$A_{ij} = \begin{cases} 0, & \text{no communication between agents } ij \\ 1, & \text{communication between agents } ij \end{cases}$$

repeat

 Create $mat_{his} = \text{concat}(X_{his}, p)$.

 Add S to F_D by (6) and (7).

 Create Factor Graph network $G = (V, F_i, F_D, A)$.

for $i = N, \dots, 1$ **do**

 Get c_i from V .

 Get $c_{F_i}^m$ from F_i .

 Get c_{F_D} from F_D .

 Create Q, K, V from (8), (9) and (10).

 Calculates the agent information.

 Get Y_{pred} by (11)

end for

until $\text{argmin}(L) = \text{RMSE}(|Y^T - Y_{pred}|^2)$

the motion trend at time $t + 1$, it is necessary to take the observation state at the current time as the query object and search for similar motion trend in the state at the historical time. Therefore, this paper takes the input characteristics and mutual information of the current moment V, F_i as the query value, and the dynamic factor of the historical moment F_D as the index and the queried value for self-attention calculation, as shown in Figure 6. For the characteristic information c_i of variable node i and the internal factor nodes $c_{F_i} = \{c_{F_i}^1, c_{F_i}^2, \dots, c_{F_i}^{m_i}\}$ of self-information connected to it, m_i is the number of factor nodes connected to variable node i . Calculate the query value of the input feature at time $t + 1$, and

obtain the query value through matrix splicing and multi-layer neural network:

$$Q = W^{(1)}(\omega_Q, W^{(2)}(\text{cat}(c_i, c_{F_i}^1, c_{F_i}^2, \dots, c_{F_i}^{m_i}), \omega_Q)) \quad (8)$$

where $W^{(1)}, W^{(2)}$ is the trainable network, ω_Q is the network parameter, and cat is the eigenmatrix concat function. For the dynamic factor node feature c_{F_D} of the hidden state S of historical motion, as the queried value at time $t + 1$, the index and value of self-attention network are calculated respectively:

$$K = W^{(3)}(\omega_K, c_{F_D}) \quad (9)$$

$$V = W^{(4)}(\omega_V, c_{F_D}) \quad (10)$$

Based on the obtained values of Q, K and V , self-attention network calculation is carried out to obtain the position feature information at time $t + 1$:

$$Y_{pred} = \text{MLP}(W^{(5)}(\text{softmax}(\frac{Q^T K}{\sqrt{d}})V), \omega) \quad (11)$$

where $W^{(5)}$ is trainable network, ω is network parameter, MLP is fully connected multilayer perceptron. In order to optimize the network weight W , root mean square error (RMSE) can be minimized. The constrained optimization objective of the loss function is:

$$\text{argmin}(L) = \text{RMSE}(|Y - Y_{pred}|^2) \quad (12)$$

In this paper, the ADAM optimizer is used to train the network model. In the training process, the mean square error between the model's prediction results of the target location and the actual location information is minimized on the training set. Finally, the prediction is made on the test set, and the positioning result of the target position is obtained.

IV. EXPERIMENT AND ANALYSIS

In this section, we use a data set of simulation scenarios to evaluate the performance of the model in this article and experiment with the following research questions (RQs).

RQ 1: What scale of multi-agent network does this model apply to, and how does the model performance change with the number of agents? (Figure 7)

RQ 2: Does the hidden state extracted from the time series of the model in this paper affect the final position estimation accuracy of the model, and how to choose the length of the time window? (Figure 8)

RQ 3: How does the communication threshold between agents affect the accuracy of the model? Under what communication conditions can the proposed model achieve optimal performance? (Figure 9)

RQ 4: What is the impact of each data module, including data fusion, internal factor and dynamic factor, on the accuracy of the model? (Figure 10)

RQ 5: Compared with existing collaborative network localization methods, how does the proposed model perform? (Figure 12)

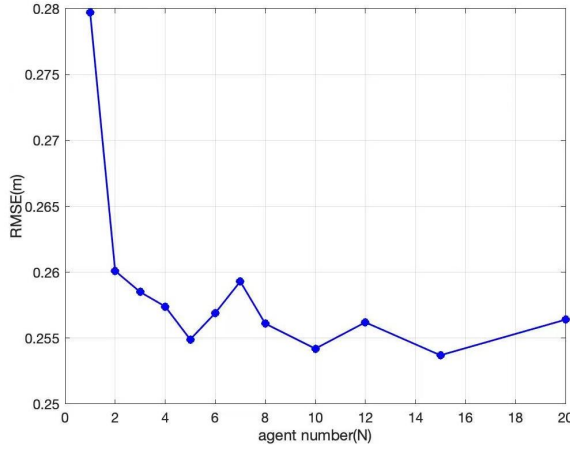


Fig. 7: Model accuracy under different number of agents.

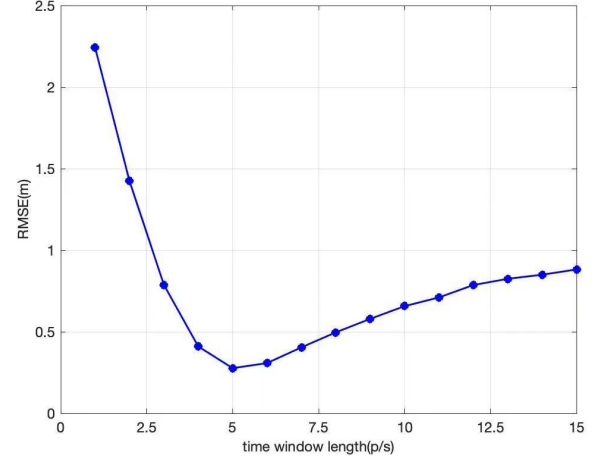


Fig. 8: Model experiment results under different window sizes.

A. Setups

Prior to initiating model training, we established a well-defined training procedure. On the hardware front, all simulations were executed on a server computer furnished with a 13-core 3.40GHz CPU and an NVIDIA 4080 GPU with 12GB of memory. Regarding the software environment, we employed a training regimen of 100 epochs for all methods under consideration, with a batch size of 15 for each epoch. The optimization was performed using the Adam optimizer, set at a learning rate of 0.001. For the parameter experiments, we selected an array of continuous parameter values to evaluate the average performance of the trained model across 100 generated test datasets. In the ablation and comparative experiments, we applied the optimal parameter values derived from the initial parameter experiments to facilitate a fair comparison across various noise conditions.

We generated the data set through the simulation scene. Firstly, the initial coordinates of N agents were randomly generated in a square area of 100m * 100m. The value of N was set according to different experimental parameters, ranging from 1-20. For each agent i , we set a fixed step size a_i and a fixed Angle b_i , and add measurement errors when it is used as sensor data, which conform to the Gaussian distribution. where, $a_i^{nois} = a_i + n_{a_i}$, $b_i^{nois} = b_i + n_{b_i}$, $n_{a_i} \sim N(0, \sigma_{a_i}^2)$, $n_{b_i} \sim N(0, \sigma_{b_i}^2)$.

In the process of data set generation, we asked all agents to walk 25 time steps in the simulation scene according to their fixed step size and angle. At each time step, we calculate the distance dis_{ij}^t between the agents as one of the sensor data, and add the measurement error $dis_nois_{ij}^t = dis_{ij}^t + n_{dis_{ij}^t}$, $n_{dis_{ij}^t} \sim N(0, \sigma_{dis_{ij}^t}^2)$ when adding the sensor data.

B. Numerical Simulations

1) *With different numbers of agents:* Aiming at problem 1 in the experimental part, we tested the positioning accuracy changes of the proposed model under different number of agents. Experimental results are shown in Figure 7. Root mean square error (RMSE) is used as the measurement standard of positioning accuracy. Our experimental analysis indicates that the model's positioning accuracy becomes notably stable

with an average error of 0.2561m when the agent count is above eight. This stability continues with increasing agent numbers, suggesting that the interaction among a larger set of agent features is beneficial up to a point. The chosen interval of [4,7] for the historical moment window size stems from a comprehensive evaluation of the model's performance across various configurations. It represents a pragmatic range within which the factor graph network can effectively utilize collaborative information to discern the agents' motion states without incurring the drawbacks of processing a potentially overwhelming feature set.

In conclusion, the model presented in this paper can show stable average performance with the increase of the number of agents in the multi-agent co-operating scenario. Compared with the traditional static positioning method, the proposed model can realize the dynamic positioning of a certain scale of agents, and can meet the actual requirements of high dynamic collaborative positioning scenarios.

2) *Effect of Different Window Sizes:* Figure 8 shows the positioning performance of this model based on historical time series data under different window sizes. We use root mean square error (RMSE) as a measure of positioning accuracy. The experimental results show that when the length of the historical time window is less than 2, the positioning accuracy of the model is the lowest, reaching 1.424 meters. With the increase of window length, the positioning accuracy of the model is gradually improved, and the minimum error of 0.2768 meters is reached when the window length is 5. However, as the window size further increases, the performance of the model begins to slowly decline. When the window length is less than 2, due to the small amount of historical data, the model cannot effectively obtain the historical motion state, resulting in poor performance. With the increase of window length, the historical state that can be obtained by the model increases, which improves the prediction effect of the motion trend of the agent. However, when the window length is further increased, the experimental results show that the more distant historical data has a smaller effect on the agent position prediction at the current moment. In addition, the number of input historical features is large, which makes it difficult for

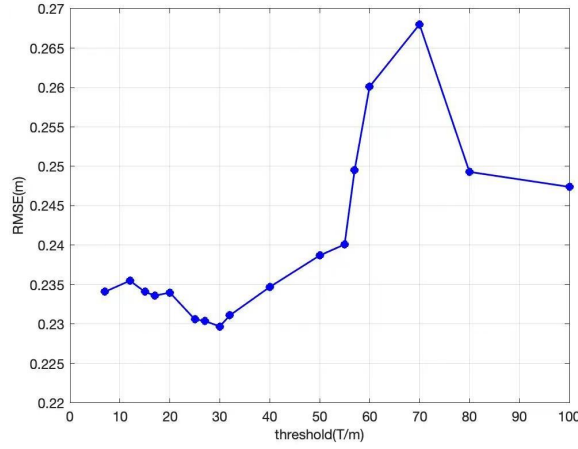


Fig. 9: RMSE under different communication thresholds.

the model to accurately extract the history hidden state from the large historical data, resulting in lower test accuracy.

In conclusion, the proposed model can effectively utilize historical data information and improve the positioning accuracy of collaborative agents by extracting the hidden state of historical motion. In terms of the selection of time window length, a reasonable window length value should be selected in the interval [4,7] according to the experimental results.

3) *Effect of Different Communication Thresholds:* In the location scenario, different communication thresholds T play an important role in data fusion and filtering. The main research in this section is to evaluate the positioning accuracy of this model under different threshold selection conditions. When the distance between two agents is greater than T , the agents do not communicate with each other. When the distance between agents is less than or equal to T , communication is connected. Figure 9 shows the results of mean square error RMSE under different values of T in this model. It can be seen that when the communication threshold T is less than 20, the mean square error RMSE is 0.2341 m. As T gradually increases to 30, RMSE decreases to the minimum value 0.2297m. When the threshold value T is further increased to greater than 55, the average error of the model reaches the maximum value of 0.2524m.

Our analysis has revealed that when the communication threshold window length is less than 25, the model is constrained by an insufficient historical data pool, which hampers its performance. As the window length extends beyond this lower threshold, the model gains access to an increasing amount of historical state information, thereby enhancing its ability to predict the agents' motion trajectory more effectively. Conversely, we observed that expanding the window length beyond 32 yields diminishing returns. This is due to the fact that older historical data exerts a progressively minor influence on the current position prediction of the agent. Additionally, a larger set of input historical features could overwhelm the model, complicating its task of accurately extracting the relevant historical states from a vast data trove, which our tests showed can lead to a decline in accuracy.

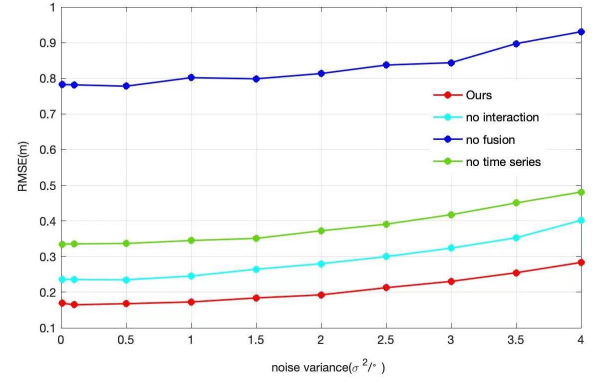


Fig. 10: RMSE under different noise variances.

TABLE II: The averaged loss (RMSE) of all methods under different angle noise conditions.

Models/angle noise(σ^2)	0.5	1	1.5	2	4
GCN	1.0346	1.1644	1.1076	1.0678	1.1294
GAT	0.4838	0.4870	0.5089	0.5290	0.5279
LS	0.5158	0.5202	0.5428	0.5441	0.6113
GTN	0.4103	0.4097	0.4172	0.4443	0.5491
Ours	0.1709	0.1785	0.1919	0.2118	0.3257

It is this nuanced relationship between window length and historical data utility that guided our selection of the interval [25,32]. This range represents the empirical sweet spot where the model benefits from sufficient historical information to enhance predictive accuracy without being encumbered by an overload of less influential data points.

4) *Ablation Study:* In order to investigate the influence of different modules in the model in this paper on the results, we conducted ablation experiments. We conducted experiments on data fusion module, historical time series data module and interaction data module between agents under different noise variance conditions, and the results are shown in Figure 10. With the increase of noise variance, the RMSE of all methods increases gradually. The model in this paper has the best performance in terms of positioning accuracy, and can also show better performance when the noise variance is high. In contrast, the model without the information fusion module performed the worst. The results of ablation experiment show that the data fusion module has the greatest influence on the model. Considering that the heterogeneous nature of the input data leads to different orders of magnitude and units of data, the convergence effect of the model is poor without data fusion. The historical data module and the interagent data module also greatly improve the model performance. When there is no historical data, the model only uses the interactive information between the current motion model and the agent to predict, resulting in poor prediction effect. However, when there is no inter-agent cooperative information, the model only uses the historical motion trend and current motion state to predict, without inter-agent information cooperative correction, the effect is also inferior to the model in this paper.

TABLE III: Multi-dimensional comparison of different methods.

Models	MSE (m)	Mean (m)	Computational time (s)
GCN	1.094	1.1046	29.214
GAT	0.5073	0.6540	27.344
LS	0.5468	0.5428	35.238
GTN	0.4461	0.5180	16.1215
Ours	0.2157	0.3376	16.734

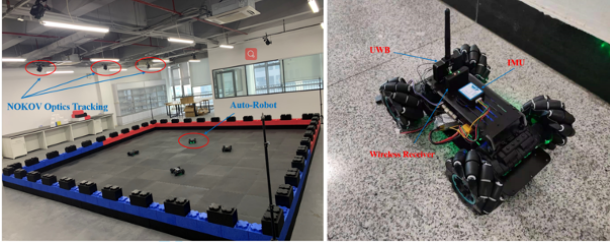


Fig. 11: (a) The scene for physical experiment. (b) The auto-robot and sensors used for the test.

C. Physical Experiment

In the physical experiment, we set up four auto-robots equipped with UWB and IMU sensors in a 10m x 10m field, as shown in Fig. 11. The robots were allowed to move in random walk motion, and their sensor information and position data were collected during the experiments. The ground truth positions were captured using the NOKOV optics motion tracking system [50]. We evaluated the algorithm's performance based on the trajectory paths and localization errors of the agents.

To assess the effectiveness of our proposed model, we conducted comparative analyses with several collaborative localization methods, including Graph Attention Networks (GAT), Least Squares (LS), Graph Convolutional Networks (GCN), and the recently introduced Graph Transformer Networks (GTN). Particularly, we utilized the GTN method as a comparative benchmark for the self-attention component within our model. In scenarios with varying noise variances and a fixed agent count of 10, we evaluated the performance of these methods, as depicted in the revised Figure 12. Our method consistently exhibited superior Mean Square Error (MSE) performance, especially with increasing noise variance. Notably, our model, along with the GTN and GAT methods, displayed a similar error trend. In contrast, the LS method showed a more gradual error increase with rising noise variance, while the GCN method demonstrated the least effective performance.

Table III presents a detailed comparison of our method against existing methods, considering metrics such as MSE, Mean Error, and computational time. The data in Table III reveals that our model achieves lower MSE and Mean Error rates, thus surpassing both the GAT and GTN methods in terms of accuracy. We also assessed computational efficiency by measuring the model training time. The results confirmed that our self-attention mechanism reduces computational time, enabling our model to function more efficiently compared to

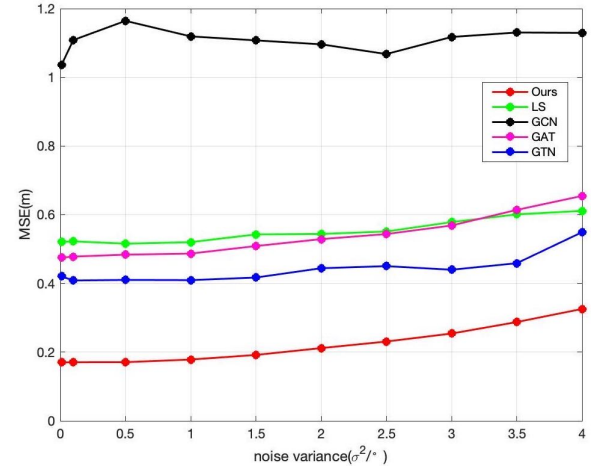


Fig. 12: Under different communication threshold, model experiment results

other attention and graph-based methods. This efficiency leads to a shorter startup time, enhancing the practical applicability of our model.

The experimental comparisons address Research Question 5 (RQ 5), illustrating that our model, along with GTN and GAT, follows a similar error trajectory. The LS method exhibits a steadier error increase as noise variance intensifies, while the GCN method is the least efficient. Importantly, our model outperforms the GTN method, which also incorporates an attention mechanism, in positioning accuracy. These findings validate that our self-attention method excels in data filtering, and our integration of this method within a factor graph-based neural network framework significantly improves localization accuracy. Compared to the static graph network approach of GCN, our model effectively utilizes historical dynamics to discern the hidden motion states of agents, facilitating more efficient data fusion and reasoning via the factor graph neural network in dynamic environments. Furthermore, our model's comparison with GAT highlights the efficacy of our self-attention method in capturing historical motion states and its dynamic adaptability to varying agent numbers, meeting the demands of dynamic collaborative scenarios.

V. CONCLUSION

In this paper, we introduce a pioneering self-attention co-location method based on factor graphs, specifically engineered to integrate multi-source heterogeneous data within a cooperative network. This approach utilizes a self-attention neural network to discern and analyze the inference processes among various types of information and actively filter agent data. Through an array of experiments, including parametric, ablation, and comparative studies, our model has proven to outperform traditional methods like Graph Convolutional Networks (GCN), Graph Attention Networks (GAT), and Multi-Layer Perceptrons (MLP). The results from these experiments underscore the influence of varying window sizes and communication thresholds on the model's performance, highlighting the strengths of our proposed approach.

However, our experiments have also uncovered certain limitations in the model. A notable issue is the model's reduced stability in accuracy when dealing with a small number of agents or when constrained by short startup times. This observation points to a need for enhancements in the model's generalization capabilities to ensure its effectiveness in a wider array of scenarios.

Moving forward, our future work will focus on overcoming this limitation. We plan to investigate strategies to boost the model's adaptability, especially in scenarios involving fewer agents and reduced operational durations. Possible solutions might include incorporating more sophisticated machine learning techniques for improved feature extraction and developing advanced training protocols to increase the model's robustness. In addition, we aim to explore the scalability of our model across diverse operational environments, with the goal of expanding its practical applicability and overall efficacy.

REFERENCES

- [1] X. Li and F. Cao, "Location Based TOA Algorithm for UWB Wireless Body Area Networks," in *2014 IEEE 12th International Conference on Dependable, Autonomic and Secure Computing*. Dalian, China: IEEE, Aug. 2014, pp. 507–511. [Online]. Available: <http://ieeexplore.ieee.org/document/6945742/>
- [2] A. Chehri, P. Fortier, and P.-M. Tardif, "On the TOA Estimation for UWB Ranging in Complex Confined Area," in *2007 International Symposium on Signals, Systems and Electronics*. Montreal, QC, Canada: IEEE, Jul. 2007, pp. 533–536. [Online]. Available: <http://ieeexplore.ieee.org/document/4294530/>
- [3] C. ZHANG, W. WANG, C. XU, X. SUN, and J. GUO, "A TOA-based optimization positioning algorithm for nonlinear sight errors," *Journal of Nanjing University of Posts and Telecommunications (Natural Science Edition)*, vol. 42, no. 4, pp. 56–63, 2022.
- [4] W. HU and Y. ZHOU, "Research on indoor localization algorithm based on UWB and IMU information fusion," no. 02 vo 45, pp. 193–197+213, 2023.
- [5] Z. Wang, Y. Wu, and Q. Niu, "Multi-Sensor Fusion in Automated Driving: A Survey," *IEEE Access*, vol. 8, pp. 2847–2868, 2020. [Online]. Available: <https://ieeexplore.ieee.org/document/8943388/>
- [6] H. Wymeersch, J. Lien, and M. Z. Win, "Cooperative Localization in Wireless Networks," *Proceedings of the IEEE*, vol. 97, no. 2, pp. 427–450, Feb. 2009. [Online]. Available: <http://ieeexplore.ieee.org/document/4802193/>
- [7] W. Yan, D. Jin, Z. Lin, and F. Yin, "Graph Neural Network for Large-Scale Network Localization," in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. Toronto, ON, Canada: IEEE, Jun. 2021, pp. 5250–5254. [Online]. Available: <https://ieeexplore.ieee.org/document/9414520/>
- [8] T. N. Kipf and M. Welling, "Semi-Supervised Classification with Graph Convolutional Networks," Feb. 2017, arXiv:1609.02907 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1609.02907>
- [9] W. L. Hamilton, R. Ying, and J. Leskovec, "Inductive Representation Learning on Large Graphs," Sep. 2018, arXiv:1706.02216 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1706.02216>
- [10] P. Velickovi, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph Attention Networks," in *6th ICLR 2018: Vancouver, BC, Canada*, Feb. 2018.
- [11] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How Powerful are Graph Neural Networks?" Feb. 2019, arXiv:1810.00826 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1810.00826>
- [12] H. Loeliger, "An Introduction to factor graphs," *IEEE Signal Processing Magazine*, vol. 21, no. 1, pp. 28–41, Jan. 2004. [Online]. Available: <http://ieeexplore.ieee.org/document/1267047/>
- [13] P. Mirowski and Y. LeCun, "Dynamic Factor Graphs for Time Series Modeling," in *Machine Learning and Knowledge Discovery in Databases*, W. Buntine, M. Grobelnik, D. Mladenić, and J. Shawe-Taylor, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, vol. 5782, pp. 128–143, series Title: Lecture Notes in Computer Science. [Online]. Available: http://link.springer.com/10.1007/9783642041747_9
- [14] Y. Yuan, X. Weng, Y. Ou, and K. Kitani, "AgentFormer: Agent-Aware Transformers for Socio-Temporal Multi-Agent Forecasting," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. Montreal, QC, Canada: IEEE, Oct. 2021, pp. 9793–9803. [Online]. Available: <https://ieeexplore.ieee.org/document/9710708/>
- [15] D. Koller and N. Friedman, *Probabilistic graphical models: principles and techniques*, ser. Adaptive computation and machine learning. Cambridge, MA: MIT Press, 2009.
- [16] Z. Zhang, F. Wu, and W. S. Lee, "Factor Graph Neural Network," Jun. 2019, arXiv:1906.00554 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1906.00554>
- [17] M. Gori, G. Monfardini, and F. Scarselli, "A new model for learning in graph domains," in *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, vol. 2. Montreal, Que., Canada: IEEE, 2005, pp. 729–734. [Online]. Available: <http://ieeexplore.ieee.org/document/1555942/>
- [18] D. Bahdanau, K. Cho, and Y. Bengio, "Neural Machine Translation by Jointly Learning to Align and Translate," May 2016, arXiv:1409.0473 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1409.0473>
- [19] Y. Kim, C. Denton, L. Hoang, and A. M. Rush, "Structured Attention Networks," Feb. 2017, arXiv:1702.00887 [cs]. [Online]. Available: <http://arxiv.org/abs/1702.00887>
- [20] A. Vaswani, "Attention Is All You Need," in *Advances in neural information processing systems 30: 31st Annual Conference on Neural Information Processing Systems (NIPS 2017): Long Beach, California, USA, 4-9 December 2017*, Dec. 2017, pp. 6000–6010.
- [21] N. Patwari, A. Hero, M. Perkins, N. Correal, and R. O'Dea, "Relative location estimation in wireless sensor networks," *IEEE Transactions on Signal Processing*, vol. 51, no. 8, pp. 2137–2148, Aug. 2003. [Online]. Available: <http://ieeexplore.ieee.org/document/1212671/>
- [22] A. Simonetto and G. Leus, "Distributed Maximum Likelihood Sensor Network Localization," *IEEE Transactions on Signal Processing*, vol. 62, no. 6, pp. 1424–1437, Mar. 2014. [Online]. Available: <http://ieeexplore.ieee.org/document/6725647/>
- [23] J. A. Costa, N. Patwari, and A. O. Hero, "Distributed weighted-multidimensional scaling for node localization in sensor networks," *ACM Transactions on Sensor Networks*, vol. 2, no. 1, pp. 39–64, Feb. 2006. [Online]. Available: <https://dl.acm.org/doi/10.1145/1138127.1138129>
- [24] P. Biswas, T.-C. Lian, T.-C. Wang, and Y. Ye, "Semidefinite programming based algorithms for sensor network localization," *ACM Transactions on Sensor Networks*, vol. 2, no. 2, pp. 188–220, May 2006. [Online]. Available: <https://dl.acm.org/doi/10.1145/1149283.1149286>
- [25] P. Tseng, "SecondOrder Cone Programming Relaxation of Sensor Network Localization," *SIAM Journal on Optimization*, vol. 18, no. 1, pp. 156–185, Jan. 2007. [Online]. Available: <http://epubs.siam.org/doi/10.1137/050640308>
- [26] A. Ihler, J. Fisher, R. Moses, and A. Willsky, "Nonparametric belief propagation for self-localization of sensor networks," *IEEE Journal on Selected Areas in Communications*, vol. 23, no. 4, pp. 809–819, Apr. 2005. [Online]. Available: <http://ieeexplore.ieee.org/document/1413473/>
- [27] D. Jin, F. Yin, C. Fritsche, F. Gustafsson, and A. M. Zoubir, "Bayesian Cooperative Localization Using Received Signal Strength With Unknown Path Loss Exponent: Message Passing Approaches," *IEEE Transactions on Signal Processing*, vol. 68, pp. 1120–1135, 2020. [Online]. Available: <https://ieeexplore.ieee.org/document/8974236/>
- [28] F. Yin, C. Fritsche, D. Jin, F. Gustafsson, and A. M. Zoubir, "Cooperative Localization in WSNs Using Gaussian Mixture Modeling: Distributed ECM Algorithms," *IEEE Transactions on Signal Processing*, vol. 63, no. 6, pp. 1448–1463, Mar. 2015. [Online]. Available: <http://ieeexplore.ieee.org/document/7015606/>
- [29] H. Chen, G. Wang, Z. Wang, H. C. So, and H. V. Poor, "Non-Line-of-Sight Node Localization Based on Semi-Definite Programming in Wireless Sensor Networks," *IEEE Transactions on Wireless Communications*, vol. 11, no. 1, pp. 108–116, Jan. 2012. [Online]. Available: <http://ieeexplore.ieee.org/document/6087384/>
- [30] H. XIE and H. YANG, "An indoor location method combining random forest with deep learning," *Journal of Shanghai Maritime University*, no. 03 vo 41, pp. 117–121, 2020.
- [31] Z. Liu, C. Huang, Y. Yu, P. Song, B. Fan, and J. Dong, "Dynamic Representation Learning for Large-Scale Attributed Networks," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. Virtual Event Ireland: ACM, Oct. 2020, pp. 1005–1014. [Online]. Available: <https://dl.acm.org/doi/10.1145/3340531.3411945>
- [32] X. Wang, H. Ji, C. Shi, B. Wang, Y. Ye, P. Cui, and P. S. Yu, "Heterogeneous Graph Attention Network," in *The World Wide Web*

- Conference*. San Francisco CA USA: ACM, May 2019, pp. 2022–2032. [Online]. Available: <https://dl.acm.org/doi/10.1145/3308558.3313562>
- [33] A. Patwardhan, R. Murai, and A. J. Davison, “Distributing Collaborative Multi-Robot Planning with Gaussian Belief Propagation,” *IEEE Robotics and Automation Letters*, vol. 8, no. 2, pp. 552–559, Feb. 2023, arXiv:2203.11618 [cs]. [Online]. Available: <http://arxiv.org/abs/2203.11618>
 - [34] R. Murai, J. Ortiz, S. Saeedi, P. H. J. Kelly, and A. J. Davison, “A Robot Web for Distributed Many-Device Localization,” *IEEE Transactions on Robotics*, vol. 40, pp. 121–138, 2024. [Online]. Available: <https://ieeexplore.ieee.org/document/10286058/>
 - [35] “Low-Density Parity-Check Codes,” in *Fundamentals of Codes, Graphs, and Iterative Decoding*. Boston: Kluwer Academic Publishers, 2002, vol. 714, pp. 137–175, series Title: The International Series in Engineering and Computer Science. [Online]. Available: http://link.springer.com/10.1007/0306477947_8
 - [36] F. Kschischang, B. Frey, and H.-A. Loeliger, “Factor graphs and the sum-product algorithm,” *IEEE Transactions on Information Theory*, vol. 47, no. 2, pp. 498–519, Feb. 2001. [Online]. Available: <http://ieeexplore.ieee.org/document/910572/>
 - [37] P. Chavali and A. Nehorai, “Distributed Power System State Estimation Using Factor Graphs,” *IEEE Transactions on Signal Processing*, vol. 63, no. 11, pp. 2864–2876, Jun. 2015. [Online]. Available: <http://ieeexplore.ieee.org/document/7060661/>
 - [38] R. Tanner, “A recursive approach to low complexity codes,” *IEEE Transactions on Information Theory*, vol. 27, no. 5, pp. 533–547, Sep. 1981. [Online]. Available: <http://ieeexplore.ieee.org/document/1056404/>
 - [39] L. Von Stumberg, V. Usenko, and D. Cremers, “Direct Sparse Visual-Inertial Odometry Using Dynamic Marginalization,” in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. Brisbane, QLD: IEEE, May 2018, pp. 2510–2517. [Online]. Available: <https://ieeexplore.ieee.org/document/8462905/>
 - [40] V. G. Satorras and M. Welling, “Neural Enhanced Belief Propagation on Factor Graphs,” Mar. 2021, arXiv:2003.01998 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/2003.01998>
 - [41] Y. Shi, B. Wang, Y. Yu, X. Tang, C. Huang, and J. Dong, “Robust anomaly detection for multivariate time series through temporal gcn and attention-based vae,” *Knowledge-Based Systems*, p. 110725, 2023.
 - [42] A. Parikh, O. Täckström, D. Das, and J. Uszkoreit, “A Decomposable Attention Model for Natural Language Inference,” in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Austin, Texas: Association for Computational Linguistics, 2016, pp. 2249–2255. [Online]. Available: <http://aclweb.org/anthology/D16-1244>
 - [43] R. Paulus, C. Xiong, and R. Socher, “A Deep Reinforced Model for Abstractive Summarization,” Nov. 2017, arXiv:1705.04304 [cs]. [Online]. Available: <http://arxiv.org/abs/1705.04304>
 - [44] S. Li, W. Chen, B. Yan, Z. Li, S. Zhu, and Y. Yu, “Self-supervised contrastive representation learning for large-scale trajectories,” *Future Generation Computer Systems*, 2023.
 - [45] V. Kosaraju, A. Sadeghian, R. Martín-Martín, I. Reid, S. H. Rezatofighi, and S. Savarese, “Social-BiGAT: Multimodal Trajectory Forecasting using Bicycle-GAN and Graph Attention Networks,” Jul. 2019, arXiv:1907.03395 [cs]. [Online]. Available: <http://arxiv.org/abs/1907.03395>
 - [46] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, “Social LSTM: Human Trajectory Prediction in Crowded Spaces,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Las Vegas, NV, USA: IEEE, Jun. 2016, pp. 961–971. [Online]. Available: <http://ieeexplore.ieee.org/document/7780479/>
 - [47] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi, “Social GAN: Socially Acceptable Trajectories with Generative Adversarial Networks,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. Salt Lake City, UT: IEEE, Jun. 2018, pp. 2255–2264. [Online]. Available: <https://ieeexplore.ieee.org/document/8578338/>
 - [48] S. Hochreiter and J. Schmidhuber, “Long Short-Term Memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997. [Online]. Available: <https://direct.mit.edu/neco/article/9/8/1735-1780/6109>
 - [49] I. Schwartz, S. Yu, T. Hazan, and A. G. Schwing, “Factor Graph Attention,” in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Long Beach, CA, USA: IEEE, Jun. 2019, pp. 2039–2048. [Online]. Available: <https://ieeexplore.ieee.org/document/8953801/>
 - [50] “Nokov,” <https://en.nokov.com/direct>, accessed on 2023-07-21.

Cheng Xu received the B.E., M.S. and Ph.D. degree from the University of Science and Technology Beijing (USTB), China in 2012, 2015 and 2019 respectively. He is currently working as an associate professor at University of Science and Technology Beijing. He is supported by the Post-doctoral Innovative Talent Support Program from Chinese government in 2019. He is an associate editor of International Journal of Wireless Information Networks. His research interests now include swarm intelligence and multi-robots network.

Ran Su is currently working towards his Master degree in University of Science and Technology Beijing (USTB), China. His research interests include swarm intelligence and internet of things.

Ran Wang received the B.E. degree from the Beijing Information Science and Technology University, China in 2013, and the M.S. degree from the University of Science and Technology Beijing (USTB), China in 2016. She is currently working toward the Doctoral degree at University of Science and Technology Beijing. Her research interests include swarm intelligence, distributed security and internet of things.

Shihong Duan received Ph.D. degree in computer science from University of Science and Technology Beijing (USTB), in 2012. She is an associate professor with the School of Computer and Communication Engineering, USTB. Her research interests include wireless indoor positioning, swarm intelligence and internet of things.