

# **Deep AI military staf: cooperative battlefeld situation awareness for commander's decision making**

**Chang‑Eun Lee<sup>1</sup> · Jaeuk Baek1 · Jeany Son<sup>2</sup> · Young‑Guk Ha3**

Accepted: 5 October 2022 / Published online: 31 October 2022 © The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2022

# **Abstract**

There are many studies adopting artifcial intelligence (AI) to develop core technologies for the future army but they are still at the level of basic research. It is expected that military power will be negatively afected by aging and declining population. In addition, as more than 500,000 agents will be dispatched to monitor combat scenes, the data sensed by each agent should be managed simultaneously recognize and evaluate the situation on the battlefeld in real time. Despite increased complexity in the battlefeld, current command system entirely rely on the experience and expertise of individual commanders, which severely restricts defense capabilities. Therefore, AI based military staff needs to be developed to identify potential threats that commanders are likely to miss, to develop smart command systems, and to provide data-driven rationale for commander's decisions. In this paper, we propose a deep AI military staff to support commander decision-making. Our proposed model is composed of four key parts: multi-agent based manned-unmanned collaboration architecture (MACA), robust tactical map fusion technology in poor environments (RTMF), hypergraph based representation learning (HRL) and space-time multi layer model for battlefelds recognition (STBR). We design an architecture and generate dataset for training the core network. Simulation results are provided to demonstrate the performance of Deep AI military staf.

**Keywords** Cooperative battlefeld situation awareness · Hidden enemy visualization · Fast panoptic segmentation · Hypergraph · Commander's decision making · AI military staf

 $\boxtimes$  Young-Guk Ha ygha@konkuk.ac.kr

Extended author information available on the last page of the article

#### **1 Introduction**

In the battlefeld environments where global navigation satellite system (GNSS) is not applicable, combatants encounter many buildings and obstacles that are not recognized beforehand. To obtain spatial information, many sensors with diverse capabilities are used by agents, but the data sufers from quality degradation due to irregular and dynamic motions of combatants. Even worse, more than 500,000 agents will be dispatched in the near future, so the data sensed by each agent should be managed simultaneously to recognize the battle situation in real time manner. In such complex battlefeld environments, it is necessary to develop manned/unmanned collaboration system for commander's decision making. In this paper, we propose a deep AI military staf for manned/unmanned agent collaboration system. The proposed deep AI military staff can create spatial map from visual and location information between agents and target location. Also, it can analyze threat in operation areas and alleviate cognitive burden of commanders by data-driven autonomous decision. To recognize global situation of the entire combat scenes, we adopt multitask scheduling with re-planning, which verifes whether mission of each agent is successful under the cyclic operation structure.

The proposed deep AI military staff is composed of four key parts: multi-agent based manned-unmanned collaboration architecture (MACA), robust tactical map fusion technology in poor environments (RTMF), hypergraph based representation learning (HRL) and space-time multi layer model for battlefelds recognition (STBR). To develop each part of our proposed deep AI military staf, we design an architecture and generate a dataset for training the core network. Simulation results are provided to demonstrate the performance of Deep AI military staf.

#### **2 Background**

The intelligent command control system is a key power system that supports the commander's decision-making and overall management. It improves the joint chiefs of staff system, army Corps-level system, and Army division-level system by supporting timely and appropriate decision-making from various surveillance and reconnaissance assets. In 2018, the U.S. Army introduced an Integrated Visual Augmentation System (IVAS) system to enhance soldiers' awareness and signed a contract with Microsoft to supply HoloLens 2 [\[1\]](#page-26-0). IVAS aims to provide augmented reality with various information (environmental information, operational overview, surrounding terrain and building structures) necessary for soldiers to train and perform their duties, and through the AI chipset that recognizes soldiers' eyes, hands and voices by 2023.

DARPA is already developing military situational awareness technologies to support decision-making. Collection and Monitoring via Planning for Active Situational Scenarios (COMPASS) and Active Interpretation of Disparate

Alternatives (AIDA) are the most representative decision support programs. COMPASS [\[2](#page-26-1)] utilizes state-of-the-art AI technology to support the commander's judgment in consideration of complex and multilayered battlefeld situations. In addition, DARPA [[3](#page-26-2)] simultaneously analyzes unstructured data entered from various multimedia sources through the AIDA project, a decision support program using multimodal data, generating various hypotheses about events, situations, and so on to support the commander's decision making.

# **3 Related works**

## **3.1 Panoptic segmentation**

Panoptic Segmentation [\[4\]](#page-26-3) is an approach that produces both instance segmentation and semantic segmentation simultaneously. However, instance segmentation is much harder than object detection or semantic segmentation since it should capture more precise and detailed structure of the instances. In case of panoptic segmentation, which provides semantic segmentation for background as well as a object detection and seg-mentation masks for each instance, it become more difficult. A naïve approach [\[4](#page-26-3)] for panoptic segmentation is to predict instance segmentation and semantic segmentation results separately and merge those outputs into one fnal panoptic segmentation result. Recently, most approaches use a single shared Feature Pyramid Network (FPN) [\[5](#page-27-0)] as a backbone for feature extraction, and then add two branches on the feature extractor to produce final panoptic segmentation outputs  $[6-9]$  $[6-9]$  $[6-9]$ . There are two categories for instance segmentation (or detection) modules in panoptic segmentation networks; one is a two-stage approach [\[7–](#page-27-3)[9](#page-27-2)] which is motivated by Mask-RCNN [\[10](#page-27-4)], and the other is a single-stage approach  $[6, 11]$  $[6, 11]$  $[6, 11]$ , motivated by SSD  $[12]$  or YOLO  $[13]$ . In our method, we choose single-stage approach to achieve real-time speed.

## **3.2 Image completion**

Image completion, also known as image inpainting, is one of the traditional tasks in a computer vision society. It is a task of reconstructing missing regions in an image based on the overall scene. Most early approaches focus on fnding similar contents from other parts or background in the image, and copy them to fll gaps [\[14,](#page-27-8) [15\]](#page-27-9). However, these approaches fail to recover large holes since they cannot capture global context of whole scenes, and also they cannot reconstruct content not present in the background image. Recently, learning based methods [[16](#page-27-10)[–18\]](#page-27-11) are proposed to solve this problem. These method generate contents based on semantic information from a large dataset.

#### **3.3 Multi‑modal data and graph neural network**

In the future battlefeld environment where multiple manned/unmanned agents such as warrior platforms, drones, and robots are expected to sense diverse kinds of data, it is critical to process these multi-modal data (e.g., visual, voice, language, graph,

etc.) for the commander's global battlefeld environment awareness. However, many challenges have been observed because of difculties on dealing with heterogeneous data in a consistent way. Even worse, the needs to classify agents that have observed the same local situations is growing for global battlefeld situation awareness, but this is not easy because tremendous time-series data are generated by agents.

Recently, multi-modal data has been processed in a consistent way by forming and analyzing graphs and their components (i.e., nodes and their relations). Generally, it is tricky to infer information from graphs, but recent advances in graph neural networks (GNNs) enables many tasks on graphs, such as unlabeled node classifcation, edge prediction and graph classifcation. To fully utilize graph structure in node classifcation, [[19\]](#page-27-12) and [\[20](#page-27-13)] propose a random walk that generates list of nodes with similar characteristics, then learn the feature vectors of nodes and edges. [\[21](#page-27-14)] adopts a metapath to enforce random work to follow the predefned path, which outperforms the pure random walk approaches. Also, multi-head attention mechanism used in natural language process is leveraged in graph [\[22](#page-27-15)]. As for graph classifcation, [\[23](#page-27-16)] proposes hierarchical graph representation that sequentially coarsens graph upto a predefned size, then learns the feature vector of entire graph. Based on Weisfeiler-Lehman (WL) algorithm, [\[24](#page-27-17)] classifes the graph in a consistent order by sorting graph nodes with the extracted information.

#### **4 System design and major contributions**

Our proposed deep AI military staf is composed of four parts; MACA, RTMF, HRL and STBR as shown in Fig. [1](#page-4-0). MACA is the entire architecture including RTMF, HRL, and STBR. RTMF includes the ability to create battlefeld scenarios, analyze information collected by individual agents, and generate data necessary for learning. HBR includes the ability to select agents that share the same situation/environment by analyzing the similarity between graphs generated by individual agents and to discriminate and confuse the similarity between nodes in the distributed graph to create a global graph. STBR includes a technology that predicts the battlefeld situation/environment in real time by analyzing the correlation of complex knowledge distributed in time and space using a fusion knowledge graph generated by HRL. We frst focus on RTMF and HRL and provide their major contributions. Then, details of each part of deep AI military staff are illustrated in Sects. [5-](#page-6-0)[8](#page-16-0).

#### **4.1 Enemy detection and hidden enemy visualization**

Computer vision is the basis of many real-life applications. In order to deal with variety applications in the real-world, numerous computer vision algorithms, such as image classifcation [[25,](#page-27-18) [26](#page-27-19)], visual tracking [[27\]](#page-27-20), object detection [[28\]](#page-27-21) and semantic segmentation [\[29](#page-28-0)], have been developed. Thanks to recent great advance in deep learning, computer vision algorithms have shown rapid improvements in the past decade.



Multi-Agent based manned-unmanned Collaboration Architecture (MACA)

<span id="page-4-0"></span>**Fig. 1** The architecture of MACA

Those computer vision algorithms are also applied to support visual cognition of soldiers in battlefelds. In recent Integrated Visual Augmentation System (IVAS) project, Microsoft's HoloLens2 is used to support displaying important information in battlefelds, since it is crucial to understand surrounding scenes for survival of soldiers. However, scene understanding becomes more challenging due to complex, drastic changes in battlefeld environments, and target object recognition failures increase. Moreover, since most of enemies hide their bodies by cover and concealment as shown in Fig. [2,](#page-5-0) it is much hard to detect them in battle. Such frequent encounters with concealed enemies can lead to increased level of fear and tension and drop in the overall moral of the troops.

In Sect. [6,](#page-8-0) we propose a new method that exposes hidden enemies in complex scenes using panoptic segmentation and image completion methods to overcome these problems. An overall framework of our proposed RTMF is shown in Fig. [3](#page-5-1). First, to detect enemies in complex scenes, we design a real-time panoptic segmentation network which runs 33fps on 550x550 resolution without severe performance drops. The feature pyramid network (FPN) [\[5](#page-27-0)] is shared for both instance and semantic segmentation to encode shared representation for each work. We also adopt a similar network architecture to YOLACT [\[30\]](#page-28-1) which is a well-known real-time instance segmentation network, and extend it to panoptic segmentation. Second, to visualize hidden enemies, we utilize Pluralistic Image Completion network (PICNet) [[31](#page-28-2)] to reconstruct occluded body parts.

The main contributions of RTMF in Sec. [6](#page-8-0) are listed as follows:



**Fig. 2** Invisible enemies due to cover and concealment

<span id="page-5-0"></span>

<span id="page-5-1"></span>**Fig. 3** The overall framework of robust tactical map fusion technology in poor environments (RTMF). A discriminator in the image completion network is only used during training

- We propose a novel framework to reveal and visualize enemies hidden by cover and concealment.
- We implement a real-time panoptic segmentation net-work based on YOLACT to detect enemies and capture characteristics in battlefeld scenes.
- We introduce a segmentation-guided image completion method to recover occluded parts of target objects.
- We demonstrate that our segmentation network achieves comparable performance to the state-of-the-art methods and our image completion network is able to recover hidden enemies successfully.

## **4.2 Hypergraph for similarity analysis**

Many studies on GNNs provide successful results on their tasks, but traditional graph structured data still lacks in terms of data representation; since an edge can only connect at most two nodes, it cannot fully represent a real word with many tangled nodes. To tackle these problems, the concept of hypergraph is introduced, where a hypergraph is composed of hypernodes and hyperedges. Specifcally, a hypernode represents an entity in a hypergraph and a hyperedge denotes a set of correlated hypernodes without any restriction on their numbers and types. By

constructing a hypergraph, two or more hypernodes with diferent features can be related one another, allowing us to extract the useful relationship information. Similar to GNNs, deep learning architectures are used in hypergraph for classifcation [\[32](#page-28-3)[–34](#page-28-4)], where the embedding vector of hypernodes is learned by hyperedge convolution operation [[32\]](#page-28-3), by processing tuple of multi-modal data [\[33](#page-28-5)] and by applying multi-head attention mechanism [\[34](#page-28-4)]. However, most of these studies focus on hypernode classifcation, but do not utilize hypergraph to analyze interrelationship between distributed graphs. Since a lot of agents in battlefeld environments measure multi-modal data that can be transformed into graph, it is critical to process distributed graphs by integrating all geographically similar information for global situation awareness.

In Sect. [7,](#page-10-0) we propose a HRL to learn the embedding vector of distributed local graphs. Based on the fact that multi-modal data sensed by agents can be transformed into graph, our objective is to classify agents that have observed the same local situations by treating each agent as a local graph and training a hypergraph-based deep learning model. After training the embedding vector of agents, we can provide adjacency matrix that shows interrelationship between agents.

The main contributions of HRL in Sect. [7](#page-10-0) are as follows.

- Based on local graph from each agent, we construct a hypergraph to integrate multiple graphs and then, utilize a hypergraph random walk to obtain a bunch of training dataset for agent embedding vector.
- Under the predefned probabilistic rules, the proposed hypergraph random walk makes random movement from one hypernode to another, which means *intergraph* (i.e. graph-to-graph) movements. Note that the conventional random walk in [\[19](#page-27-12)[–21](#page-27-14)] makes an *intra-graph* movement from one node to another in the graph.
- Because of unlabeled situation information, we train the agent embedding vector in an unsupervised manner to make the embedding vectors of neighboring agents similar. As a result of similarity analysis, we can provide adjacency matrix of agents, which can be used along with any GNNs to construct global graphs and perform graph convolution for the commander's battlefeld situation awareness.

# <span id="page-6-0"></span>**5 MACA: Multi agent‑based manned‑unmanned collaboration architecture**

Figure [4](#page-7-0) shows the overall architecture of the proposed system for enhancing awareness of combatants in a building or an underground bunker. The collaborative agent generates a collaboration plan according to a mission, requests neighboring collaborative agents to search for knowledge/devices available for collaboration and review the availability of the knowledge/devices, generates an optimal collaboration combination on the basis of a response to the request to transmit a collaboration request, and upon receiving the collaboration request, performs the mission through mutual distributed knowledge collaboration. Such a collaborative agent may provide information about systems, battlefelds, resources, and tactics through a determination



<span id="page-7-0"></span>**Fig. 4** The proposed multi-agent based manned-unmanned collaboration architecture (MACA)

intelligence processing unit, such as complicated situation recognition, coordinative simultaneous localization and mapping (C-SLAM), and a self-negotiator.

Meanwhile, in order to support a commander in command decision, the collaborative agent combines the collected 5 pieces of information to be subjected to artifcial intelligence (AI) deep learning-based global situation recognition and CSLAM technology to provide the commander with command decision information merged with unit spatial maps through the autonomous driving robot linked with the smart helmet worn by the commander.

To this end, referring to Fig. [4,](#page-7-0) the collaborative agent includes a multi-modal object data analysis unit, an inter-collaborative agent collaboration and negotiation unit, and an autonomous collaboration determination and global situation recognition unit so that the collaborative agent serves as a supervisor of the overall system.

In addition, the inter-collaborative agent collaboration and negotiation unit searches a knowledge map through a resource management and situation inference unit to determine whether a mission model that is mapped to a goal state corresponding to the situation and environment data is present, checks integrity and safety of multiple tasks in the mission, and transmits a multi-task sequence for planning an action plan for the individual tasks to an optimal action planning unit so that the tasks are analyzed and an optimum combination of devices and knowledge to perform the tasks is constructed.

On the other hand, the optimal action planning unit performs refnement/division/allocation on action-task sequences to deliver relevant tasks to the collaborative agents located in a distributed collaboration space on the basis of a generated optimum negotiation result through a knowledge/device search and connection protocol of a hyper-intelligence network formed through the autonomous driving robots so as to deliver the relevant tasks to wearers of the respective smart helmets.

In addition, the autonomous collaboration determination and global situation recognition unit verifes whether an answer for the goal state is satisfactory through global situation recognition monitoring using a delivered multitask planning sequence using a collaborative determination/inference model and, when the answer is unsatisfactory, requests the inter-collaborative agent collaboration and negotiation unit to perform mission re-planning to have a cyclic operation structure.

#### <span id="page-8-0"></span>**6 RTMF: Robust tactical map fusion technology**

This section describes the details of our proposed frame-work to visualize hidden enemies using panoptic segmentation and Segmentation-guided image completion methods. We frst detect enemies by panoptic segmentation, and then using those segmentation results as occlusion masks for the image completion network to recover and visualize hidden parts of an enemy.

#### **6.1 Real‑time panoptic segmentation**

We design a real-time panoptic segmentation network based on YOLACT [[30\]](#page-28-1), which is a well-known single-stage instance segmentation network, to detect enemies. Our panoptic segmentation network consists of three sub-networks: feature pyramid network (FPN), instance segmentation network, and a semantic segmentation network. The fnal panoptic segmentation outputs are generated by combining two outputs of instance and semantic segmentation networks, as shown in Fig. [5.](#page-8-1) Our single-stage network has the advantage of a high speed compared to the multistage network that uses the Region Proposal Network (RPN) separately. In general, there are two types of classes for panoptic segmentation; one is a Things class and the other is a Stuf class. Our instance segmentation network only produces results for a Things class, and the semantic segmentation network performs for all classes.



<span id="page-8-1"></span>**Fig. 5** Overall architecture of the proposed panoptic segmentation network. Our network consists of three sub-networks: feature pyramid network (FPN), instance segmentation network and semantic segmentation network. Final panoptic segmentation outputs are generated by merging two outputs of each segmentation network

#### **6.1.1 Instance segmentation network**

We adopt a YOLACT [\[30](#page-28-1)] architecture for out instance segmentation network. Most single-stage instance segmentation methods motivated by SSD and YOLO are fast but show worse performance than two-stage ones. To increase performance while keep real-time speed, YOLACT divides an instance segmentation problem into two parallel tasks. One (ProtoNet) is generating a set of prototype masks on images, and the other (Prediction head) is predicting a set of coefficient to compute linear combination with prototype per instance. Two separate networks for these tasks share the FPN [[5\]](#page-27-0), and we denote k-level features on FPN by Fk heads for simplicity. Object masks for the anchor boxes of each cell are generated by multiplying the prototype activation by the mask coefficient of the prediction head. Then, the final object mask is obtained by binarization through crop and threshold operations to obtain the fnal instance segmentation results.

Our loss function for training the instance segmentation network consists of three diferent losses; a bounding box regression loss, a classifcation loss, and a mask coefficient loss, as follows:

$$
L_{\text{inst}} = L_{\text{boxreg}} + L_{\text{boxcls}} + L_{\text{mask}},\tag{1}
$$

where  $L_{\text{boxreg}}$  is a smooth-L1 loss to regress offsets for box coordinates of anchors, *L*boxcls is a softmax cross-entropy loss, and *L*mask is a pixel-level binary cross-entropy loss between a linear combination of prototype masks and a ground-truth mask for each bounding box. Note that, these losses for the instance segmentation network are only computed for Things classes. For more details of the instance segmentation network, please refer to YOLACT [[30\]](#page-28-1).

#### **6.1.2 Semantic segmentation network**

To produce segmentation masks for Stuf classes, we combine a semantic segmentation branch into the panoptic segmentation network. Our semantic segmentation network is also attached to FPN and uses feature maps from F2 to F6 heads. To reduce model complexity, the semantic segmentation network employ the shared head used in the prediction head of the instance segmentation network. Thus, the shared head also used to encode semantic features for each features head. This leads to achieving real-time panoptic segmentation by avoiding heavy computation for semantic segmentation. Those features are upsampled and concatenated into one single feature map and feed to one last convolution layer to produce the fnal semantic segmentation output masks. The detailed structure of our semantic segmentation network is shown in Fig.  $6$ .

To train the semantic segmentation network, we use a cross-entropy loss with pixel-level hard-negative mining [\[35](#page-28-6)] where pixels with large losses have more penalty. In our experiments, we penalize on pixels of which losses are larger than 30% percentile of all predictions. Then, the fnal loss for training the proposed network is as follows:



<span id="page-10-1"></span>**Fig. 6** Detailed architecture of the proposed semantic segmentation network

$$
L_{tot} = L_{inst} + \lambda \cdot L_{sem},\tag{2}
$$

where  $L_{\text{sem}}$  is a pixel-level cross-entropy loss between a segmentation segmentation prediction and a ground-truth segmentation mask generated by panoptic segmentation masks, and  $\lambda$  is a constant and is set to 3 in our experiments.

#### **6.2 Segmentation‑guided image completion**

In order to visualize hidden enemies, we use an image completion network based on conditional variational auto-encoder with generative adversarial network (cVAE-GAN) which is referred to as pluralistic image completion network (PICNet) [[31\]](#page-28-2). Using the result of the proposed panoptic segmentation as an input mask, enables robust restoration against extreme occlusion. First, we fnd hidden enemies and their occluded parts using the panoptic segmentation net-work. Predicted panoptic segmentation masks on occluded parts are regarded as a mask to be recovered. This image completion network consists of a generator network that removes the hidden area and a discriminator that determines the image generated by the generator and the original image, as shown in Fig. [7](#page-11-0). After learning, only the generator is used to restore the occluded region. Please refer to PICNet [[31\]](#page-28-2) for further details.

# <span id="page-10-0"></span>**7 HRL: Hypergraph based agent representation learning**

In this section, we propose how to deal with multi-modal data in a consistent way, then construct a hypergraph for inter-graph similarity analysis. By learning the embedding vector of agent (i.e., local graph), we can obtain adjacency matrix of agents that is the key element in GNNs to perform graph convolution for the commander's battlefeld situation awareness.



<span id="page-11-0"></span>**Fig. 7** Overall system of the image completion network based on PICNet [\[31](#page-28-2)]. Occluded masks are generated by panoptic segmentation outputs

#### **7.1 Overall procedures**

Figure [8](#page-11-1) represents system architecture of our proposed hypergraph based representation learning (HRL). Based on multi-modal data sensed by agents, HRL frst converts them into graphs. Then, a hypergraph is constructed to integrate multiple graphs, and training data are made by hypergraph random walk. To train an embedding vector of agents, we minimize triplet loss, where anchors, positive samples and negative samples are selected from hypergraph induced information. After training the embedding vector of agents, adjacency matrix is made to represent interrelationship between agents, which can be used in later part of Deep AI Military Staff for the commander's battle field situation awareness. Details of our HRL are illustrated in the following subsections.



<span id="page-11-1"></span>**Fig. 8** System architecture for our hypergraph based representation learning (HRL)

# **7.2 Graph from multi‑modal data**

In this subsection, we illustrate how to convert multi-modal data into knowledge graphs. A graph is a structure of nodes and edges, where nodes with similar properties are connected by edge to represent their relational information, and a feature vector is assigned to each node to indicate its characteristics. We frst describe how to extract graph components from multi-modal data. Then, we explain the needs for a database for node feature vectors, and construct a knowledge graph.

# **7.2.1 Graph components**

From multi-modal data, we can extract entities and attributes that can be reconstructed as contextual phrases. For example, in combat scenes, some snipers with brown hair aim their guns at the enemy while other soldiers with bombs hide behind large trees. Based on an object detection algorithms [[13,](#page-27-7) [28,](#page-27-21) [37\]](#page-28-7), objects (such as snipers, soldiers, bomb, enemy and trees) can be detected in images and interpreted as nodes in a graph. In addition, visual grounding techniques can be applied to fnd attributes of each object (e.g., brown hair and large) and relationship between objects, where the former and the latter denote the attributes and edges in graph, respectively. This is an simple example of how to convert image to graph, but other types of raw data (sounds, texts, etc) can also be transformed into graph without any difficulties.

# **7.2.2 Database**

Although nodes, attributes and edges can be obtained directly from raw data, more eforts are required to obtain node feature vector (NFV) because it can be used in any other graphs to indicate unique characteristics of nodes. To this end, we construct database that saves all the entities with their feature vectors. Based on contextual phrases that can be obtained from real world (e.g., combat movies, combat games, or combat histories), we create vocabulary of entities and their feature vectors, where word embedding modules [[38\]](#page-28-8) is used to obtain the feature vector of entities (see Fig. [9](#page-13-0) for more information). Then, the trained feature vectors of entities are used as a basis to create NFVs.

# **7.2.3 Knowledge graph**

Now we are ready to construct knowledge graph from multi-modal data. Compared to conventional graphs that has nodes, attributes along with adjacency matrix among nodes, we construct knowledge graph with nodes and NFVs, where NFVs are trained to contain information about relations and attributes. Note that in battlefeld scenario where multiple distributed graphs need to be considered simultaneously, adjacency matrix of a single graph only carries local structure information and is not meaningful for the entire distributed graphs. So, we train node-relation vector to



<Contextual phrases>

<span id="page-13-0"></span>**Fig. 9** Graph components and the feature vector of entities. The image and contextual phrases are from visual genome dataset [\[36](#page-28-9)]. Many contextual phrases can be made from a single image, and each phrase has entities with attributes, which correspond to the nodes in graph. The feature vectors of entities are trained using word embedding modules and can be used to create NFVs

utilize adjacency matrix implicitly. Also, we train node-attribute vector to represent correlated nodes and attributes, which can treat nodes with and without attributes in a consistent way.

Based on feature vector of entities in database, we train node-relation and nodeattribute vectors using [[21\]](#page-27-14), and concatenate them to create NFV with fxed size.

## **7.3 Hypergraph random walk**

In this subsections, we construct a hypergraph to integrate distributed graphs and apply hypergraph random walk to obtain training data for agent embedding vector.

# **7.3.1 Hypergraph**

A hypergraph is composed of hypernodes and hyperedges, where a hypernode means an agent, and each node in distributed graphs makes a hyperedge. For example, when each of *N* agents has a graph with *M* nodes, a hypergraph is composed of *N* hypernodes and *N* × *M* hyperedges. Each hyperedge is interpreted as a set of correlated hypernodes and can be defned through similarity analysis



<span id="page-14-0"></span>**Fig. 10** Example of hypergraph

Hyperdege	<b>Hypernodes</b>	
	Related	Unrelated
	2, 4	1, 3, 5
		1, 2, 4, 5
	1, 3	2, 4, 5
		1, 2, 3, 4
	1, 2, 5	

<span id="page-14-1"></span>**Fig. 11** Example of hypergraph induced information

between nodes in distributed graphs. Specifcally, by regarding each node as a basis, NFVs from previous subsections are used to sort all nodes by similarity. Then, a threshold or a predetermined number can be used to select the nodes such that agents (i.e., hypernodes) corresponding to the selected nodes would be included in a hyperedge (see an example of hypergraph in Fig. [10](#page-14-0)).

Once a hypergraph is obtained, a hypergraph incidence matrix can be defned, where a row denotes a hypernode (i.e., an agent) and a column denotes a hyperedge (see the middle part of Fig. [8](#page-11-1)). The element of hypergraph incidence matrix can be defned as binary, where one is assigned when a hypernode is included in a hyperedge, and zero is assigned otherwise. From the hypergraph incidence matrix, we can extract relationship information among hypernodes. Figure [11](#page-14-1) represents an example of hypergraph induced information for the given hypergraph incidence matrix in Fig. [8.](#page-11-1) It is observed from Fig. [11](#page-14-1) that a single hyperdedge relates multiple hypernodes. Also, for each hyperedge, hypernodes can be divided into two sets; one for related hypernodes and others for unrelated hypernodes.

#### **7.3.2 List of hyperedges**

Based on hypergraph induced information, hypergraph random walk can be used to obtain training data for agent embedding vector. Since all nodes in distributed graphs have their own information, it is critical to visit all the nodes to carry all information about graphs. Based on that each hyperedge can be defned on each node, our proposed hypergraph random walk defnes inter-graphs movement by performing permutation of all hyperedges. Note that one permutation makes the list of hyperedges, and several independent trials of permutations provide a bunch of training data. A batch size for training a network is defned by considering all or part of the hyperedge list.

#### **7.4 Training agent embedding vector for adjacency matrix**

In this subsection, we train a network to obtain agent embedding vectors, then construct an adjacency matrix that shows relationship information between agents.

#### **7.4.1 Loss function**

The objective of training agent embedding vector is to classify agents that have observed the same local situations. We adopt triplet loss in training because minimization of this loss makes embedding vectors of two samples (i.e., anchor and positive sample) similar, while those of anchor and negative samples are trained to be dissimilar.

Let  $\Psi$  be the set of all hypernodes (i.e., agents), and  $\Psi$ <sup>*p*</sup> and  $\Psi$ <sup>*N*</sup>*N* denote positive and negative samples of hypernodes, respectively. When hypernode  $i \in \Psi$  is used as anchor, triplet loss can be expressed as

$$
L(X_i) = \sum_{p \in \Psi_P} \log \sigma \Big( f(X_i, X_p) \Big) - \sum_{n \in \Psi_N} \log \sigma \Big( f(X_i, X_n) \Big), \tag{3}
$$

where  $X_i$  is embedding vector of hypernode  $i \in \Psi$ , *f* is a vector operator for similarity analysis and  $\sigma$  is a Relu function. Note that  $\Psi_p$  and  $\Psi_p$  are updated as anchor changes, and hypergraph induced information (in Fig. [11](#page-14-1)) can be used to defne hypernode pairs of anchor, positive and negative samples. For example, for hyper-edge 1 in Figs. [10](#page-14-0) and [11,](#page-14-1) hypernode (i.e., agent) that has graph with node 1 is used as anchor, and related hypernodes (i.e., hypernodes 2 and 4) and unrelated hypernodes (i.e., hypernodes 1, 3, 5) are used as positive and negative samples, respectively.

#### **7.4.2 Adjacency matrix of agents**

Once embedding vectors of agents are obtained, we can utilize them to construct adjacency matrix, where each row and column represents the agent. Since each

agent corresponds to graph, adjacency matrix of agents can be used to construct global graph from multiple local graphs.

It is worth noting that adjacency matrix only represents relation between distributed graphs, so it seems that there is no information about which parts of adjacent graphs (i.e., nodes) are connected. However, since hypergraph induced information provides relational information in node level, we can merge multiple local graphs without difficulties.

#### <span id="page-16-0"></span>**8 STBR: Space‑time multilayer model for battlefelds recognition**

#### **8.1 Battlefeld object model**

Figure [12](#page-16-1) is a common-based ontological design for distributed agents mounted on situational cognitive systems and individual combatant equipment to identify and share situations in the battlefeld environment with the various objects present in the battlefeld environment.

The battlefeld situation is represented by the spatial interaction of various objects present in the battlefeld environment, which requires the representation of each object as a correlated vector in order to learn/recognize efectively based on a deep neural network model. In this paper, for the development of a situationaware system, we express a formal world model based on the spatial and timespace relationship between the various objects and objects that make up the battlefeld environment as shown in Fig. [13](#page-17-0).



<span id="page-16-1"></span>**Fig. 12** Battlefeld situation ontology (BSO)



**Fig. 13** World model for situation recognition

#### <span id="page-17-0"></span>**8.2 Collaborative knowledge based battlefeld situation awareness**

Existing situational awareness systems use prebuilt situational knowledge-based symbolic reasoning, making it very difficult to infer situational knowledge building or unexpected situations in complex, time-space dynamic environments such as battlefelds. in this paper, we present a new technique to continuously learn time-space changes for various situations through mechanical learning based on environmental awareness information collected from the battlefeld.

Engagement is a very important situation to be aware of in a battlefeld environment, and in order to accurately recognize the situation, it is necessary to recognize light sources such as sparks or explosions from frearms (guns, heavy weapons, etc.). Figure [14](#page-18-0) illustrates a machine learning-based process that can recognize a feld situation by recognizing the light source in real time based on the texture characteristics of various light sources that can occur in the feld.

## **9 Experiment results**

In this section, we evaluate our panoptic segmentation on the public benchmark, and compare our performance to the state-of-the-art methods. Also we show image completion results on the our private dataset constructed using the Battlefeld4.



<span id="page-18-0"></span>**Fig. 14** The machine learning-based process for situation recognition

#### **9.1 Real‑time panoptic segmentation in RTMF**

#### **9.1.1 Implementation details**

We train our model using 3 GPUs with ImageNet pretrained ResNet-50-FPN as our feature backbone, where the batch size is set to 1 per each GPU, and a input image size of  $1440 \times 720$ , We freeze BatchNorm layers in the backbone and add GroupNorm [[39\]](#page-28-10) layers after backbone. We use Adam Solver [[40](#page-28-11)] for 400k iterations with initial learning rate of 103 and weight decay of 104. Learning rate is decreasing by a factor of 0.1 at 28 and 36k.

#### **9.1.2 Datasets**

We evaluate our method on the Cityscapes panoptic segmentation benchmarks [\[41\]](#page-28-12). The Cityscapes dataset consists of street scenes with a total of 19 classes, where 8 classes for things which have instance-level labels as well as semantic class labels, and 11 classes for stuf which only have semantic labels. Following the standard protocol, we use 2975 image for training and 500 images for testing with a resolution of  $1024 \times 2048$  for evaluation. Note that we only used finegrained labels for training. For data augmentation, color jiterring, random crop, random scaling, random fipping are applied. We also test on the Battlefeld4 dataset which is constructed by images captured from the Battlefeld4 game.

<span id="page-19-0"></span>

#### <span id="page-19-1"></span>**Table 2** Comparison with the state-of-the-art methods on the cityscapes validation set



#### **9.1.3 Evaluation metrics**

To evaluate the performance for panoptic segmentation, we use the Panoptic Quality (PQ) metric [\[4](#page-26-3)] which is computed by multiplication of two sub-factors, SQ (Segmentation Quality) and RQ (Recognition Quality), as follows:

$$
PQ = \frac{\sum_{(p,q)\in TP} IoU(p,q)}{|TP|} \times \frac{|TP|}{|TP| + 0.5 \times |FP| + 0.5 \times |FN|}
$$
  
= 
$$
SQ \times RQ = \frac{\sum_{(p,q)\in TP} IoU(p,q)}{|TP| + 0.5 \times |FP| + 0.5 \times |FN|},
$$
 (4)

where *p* and *g* denotes predicted and ground-truth segmentation masks, respectively. |TP|, |FP| and |FN| denote the number of true positive, false positive and false negative, respectively. *p* is regarded as positive detection when the intersection-overunion (IoU) of *p* and *g* is greater than 0.5. We also measure instance segmentation accuracy by averaging over  $AP<sup>r</sup>$  [\[42](#page-28-13)].

#### **9.1.4 Results**

We report PQ, SQ and RQ of the propose network for all, things and stuff classes in Table [1.](#page-19-0) We also compare the performances of our method to the state-of-the-art panoptic segmentation methods in Table [2.](#page-19-1) For fair comparison, we only report the results using ResNet-50 or ResNet-50-FPN backbones. Although our method does not outperform UPSNet [[9\]](#page-27-2), our model is the fastest method among the state-of-theart methods, while the performance of our method achieves the second-best performance. In comparison with FPSNet [[6\]](#page-27-1) which shows similar speed to our method,

the performance of the proposed network gains almost 3 points in terms of a PQ metric. In Tables  $3$  and  $4$ , we show the quantitative results for Things and Stuff classes, respectively.

We demonstrate qualitative results on the Cityscapes dataset and the Battlefield 4 dataset in Figs. [15](#page-22-0) and [16,](#page-22-1) respectively. As shown in Fig. [15,](#page-22-0) our panoptic segmentation network successfully segments small and thin structures such as traffic light or traffic sign. Also even with low-quality lighting conditions, our method shows good performance on detecting and segmenting objects in simulated battlefelds.

#### **9.2 Occluded enemy reconstruction in RTMF**

#### **9.2.1 Implementation details and datasets**

We train the image completion network using the Battlefeld4 dataset which is constructed by the screenshots of the Battlefeld4 game. This dataset consists of 1300 images of various poses and situation of soldiers. All images are resized to 256x256 and uses random free-form masks for training. Image completion network is trained on a single GPU with a batch size of 20.

#### **9.2.2 Results**

In Fig. [17](#page-23-0), reconstruction results using the image completion network are illustrated, where gray regions in the masked images show examples of irregular masks used for training. Even in the case of extreme occlusion, where only a small part of the helmet is visible, the proposed method successes in reconstructing the missing parts well. We also test our method using panoptic segmentation results and show in Fig. [18](#page-24-0).

## **9.3 Agent similarity analysis in HRL**

## **9.3.1 Implementation details and datasets**

We adopt visual genome dataset  $[36]$  $[36]$  to assign a graph to each agent. To be specific, using the contextual phrases for the single image, at least one phrase is allocated to each agent, some of which has more than two phrase. Then, entities, attributes and edges are extracted from each phrase, which form a graph. Note that even in single image, many distributed graphs can be made because of multiple tangled entities.

<span id="page-20-0"></span>





SQ 97.86 97.86 92.65 92.41 76.91 76.956 92.67 76.81 76.936 92.41 92.41 92.41 92.41 RQ 99.23 98.988 98.06.88 18.999 92.71 15.986 90.1.51 11.71.05 14.71.06 20.45 95.96 93.885 98.85 93.8

<span id="page-21-0"></span>71.76<br>38.94

69.36 66.81

90.02 98.85

76.11<br>83.98



**Fig. 15** Qualitative results on the Cityscapes dataset

<span id="page-22-1"></span><span id="page-22-0"></span>

**Fig. 16** Qualitative results on the Battlefeld4 dataset. Left: input image; Right: panoptic segmentation and detection results

# Masked image Ground-truth **Generated Image**

<span id="page-23-0"></span>**Fig. 17** Reconstruction results on the Battlefeld4 dataset

#### **9.3.2 Results**

Tables [5](#page-24-1) and [6](#page-24-2) represent the pairs of agents sorted by cosine similarity of trained agent embedding vector. We considered 26 agents and assigned a graph with more than two nodes to each agent. The number of connected nodes between graphs can be a basis that represents similarity level between agents. So, we can interpret that adjacent graphs shares many nodes. It is observed from Tables [5](#page-24-1) and [6](#page-24-2) that the trained agent embedding vector has high cosine similarity for adjacent graph and low for dissimilar graphs. By setting a threshold value and classifying the agents, we can create an adjacency matrix of agents (see Fig. [19\)](#page-25-0).

Note that the adjacency matrix in Fig. [19](#page-25-0) has relation information between multiple distributed data from each agent, and can be used in our proposed space-time multilayer model for battlefelds recognition (STBR) to predict situation labels of the merged graph.



**Fig. 18** Qualitative results for reconstructing invisible parts on the Battlefeld4 dataset

<span id="page-24-1"></span><span id="page-24-0"></span>

<span id="page-24-2"></span>



<span id="page-25-0"></span>**Fig. 19** Adjacency matrix of agents. We set 0.1 as threshold of cosine similarity and adjacent agents are indicated in black in the matrix

## **9.4 Battle fled awareness in STBR**

#### **9.4.1 Implementation details and datasets**

Through the battlefeld simulation game (Battlefeld 4), object learning data was collected to build a YOLOV3-based battlefeld situation simulation object learning network. The object was recognized based on the battlefeld data from a specifc Time Step collected from four allies(agents). We integrated the recognized object information from the screen acquired from our combatants and the location information of our combatants recognized from the map screen to input the deep learning model for situational reasons. The learning data consisted of a total of 250,000 frames of image data in three battlefeld situations (Secure, Vigilant, Engagement), and a total of 1,000,000 frames of image data by using each of the four Agent PCs for one Time Step. We classify data into three battlefeld situation; a vigilant situation when enemies, tanks, and helicopters were recognized around the allies, an engagement situation when bombs and explosions were recognized, and a safety situation otherwise.

## **9.4.2 Results**

We measure the average accuracy of situation recognition using servers with the NVIDIA RTX2080. In our simulation, it is observed that average accuracy is 87.248 % for 10,000 frames with 20.45 frame per second (FPS).

# **10 Conclusion**

In this study, we have proposed a Deep AI military staff for supporting commander decision-making by enhancing awareness of combatants. We classify data into three battlefeld situation; a vigilant situation when enemies, tanks, and helicopters were recognized around the allies, an engagement situation when bombs and explosions were recognized, and a safety situation otherwise. The proposed model provides collaborative intelligence-based real-time battlefeld situation recognition technologies, which is expected to be applicable to actual battlefeld environments or combat training simulator. In addition, our proposed models can be used as key technologies not only in the defense area but also in the 4<sup>th</sup> industrial revolution such as selfdriving cars, intelligent robots, smart factory, intelligent security/crime prevention and IoT services.

**Acknowledgements** This work was supported by Electronics and Telecommunications Research Institute (ETRI) grant funded by the Korean government (21YR1610, Research on Spatio-temporal Multi-Layer Battlefield Situation Awareness for AI Military Staff) and Korea Research Institute for Defense Technology planning and advancement (KRIT) grant funded by the Korea government DAPA(Defense Acquisition Program Administration) (No. 20-107-C00-008-02, Control technology for collective operation of military ultra-small ground robots).

**Data availability** All data generated or analysed during this study are included in this published article. For more information on datasets, please see sections. 9.1.2 and 9.3.1.

# **References**

- <span id="page-26-0"></span>1. Weiss J (2020) The army's ultimate heads-up display: The ivas. [https://sofrep.com/news/the-armys](https://sofrep.com/news/the-armys-ultimateheads-up-display-the-ivas/)[ultimateheads-up-display-the-ivas/](https://sofrep.com/news/the-armys-ultimateheads-up-display-the-ivas/)
- <span id="page-26-1"></span>2. DARPA: DARPA demonstrates "Competition" tool at combatant command (2020). [https://www.](https://www.darpa.mil/news-events/2020-03-19a) [darpa.mil/news-events/2020-03-19a](https://www.darpa.mil/news-events/2020-03-19a)
- <span id="page-26-2"></span>3. DARPA: Active interpretation of disparate alternatives (AIDA) (2017). [https://www.darpa.mil/](https://www.darpa.mil/newsevents/) [newsevents/](https://www.darpa.mil/newsevents/)
- <span id="page-26-3"></span>4. Kirillov A, He K, Girshick R, Rother C, Dollár P (2019) Panoptic segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9404–9413
- <span id="page-27-0"></span>5. Lin T-Y, Dollár P, Girshick R, He K, Hariharan B, Belongie S (2017) Feature pyramid networks for object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2117–2125
- <span id="page-27-1"></span>6. de Geus D, Meletis P, Dubbelman G (2020) Fast panoptic segmentation network. IEEE Robot Autom Lett 5(2):1742–1749
- <span id="page-27-3"></span>7. Kirillov A, Girshick R, He K, Dollár P (2019) Panoptic feature pyramid networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6399–6408
- <span id="page-27-22"></span>8. Li Y, Chen X, Zhu Z, Xie L, Huang G, Du D, Wang X (2019) Attention-guided unifed network for panoptic segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7026–7035
- <span id="page-27-2"></span>9. Xiong Y, Liao R, Zhao H, Hu R, Bai M, Yumer E, Urtasun R (2019) Upsnet: A unifed panoptic segmentation network. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8818–8826
- <span id="page-27-4"></span>10. He K, Gkioxari G, Dollár P, Girshick R (2017) Mask r-cnn. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 2961–2969
- <span id="page-27-5"></span>11. Gao N, Shan Y, Wang Y, Zhao X, Yu Y, Yang M, Huang K (2019) Ssap: Single-shot instance segmentation with affinity pyramid. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 642–651
- <span id="page-27-6"></span>12. Liu W, Anguelov D, Erhan D, Szegedy C, Reed S, Fu C-Y, Berg AC (2016) Ssd: Single shot multibox detector. In: European Conference on Computer Vision, pp. 21–37. Springer
- <span id="page-27-7"></span>13. Redmon J, Divvala S, Girshick R, Farhadi A (2016) You only look once: Unifed, real-time object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 779–788
- <span id="page-27-8"></span>14. Barnes C, Shechtman E, Finkelstein A, Goldman DB (2009) Patchmatch: a randomized correspondence algorithm for structural image editing. ACM Trans Gr 28(3):24
- <span id="page-27-9"></span>15. Bertalmio M, Sapiro G, Caselles V, Ballester C (2000) Image inpainting. In: Proceedings of the 27th Annual Conference on Computer Graphics and Interactive Techniques, pp. 417–424
- <span id="page-27-10"></span>16. Satoshi I, Edgar S-S, Hiroshi I (2017) Globally and locally consistent image completion. ACM Trans Gr 36(4):3073659
- 17. Liu G, Reda FA, Shih KJ, Wang T-C, Tao A, Catanzaro B (2018) Image inpainting for irregular holes using partial convolutions. In: Proceedings of the European Conference on Computer Vision (ECCV), pp. 85–100
- <span id="page-27-11"></span>18. Pathak D, Krahenbuhl P, Donahue J, Darrell T, Efros AA (2016) Context encoders: Feature learning by inpainting. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2536–2544
- <span id="page-27-12"></span>19. Perozzi B, Al-Rfou R, Skiena S (2014) Deepwalk: Online learning of social representations. In: Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 701–710
- <span id="page-27-13"></span>20. Grover A, Leskovec J (2016) node2vec: Scalable feature learning for networks. In: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 855–864
- <span id="page-27-14"></span>21. Dong Y, Chawla NV, Swami A (2017) metapath2vec: Scalable representation learning for heterogeneous networks. In: Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 135–144
- <span id="page-27-15"></span>22. Veličković P, Cucurull G, Casanova A, Romero A, Lió P, Bengio Y (2018) Graph Attention networks
- <span id="page-27-16"></span>23. Ying R, You J, Morris C, Ren X, Hamilton WL, Leskovec J (2018) Hierarchical graph representation learning with diferentiable pooling.<http://arxiv.org/abs/1806.08804>
- <span id="page-27-17"></span>24. Zhang M, Cui Z, Neumann M, Chen Y (2018) An end-to-end deep learning architecture for graph classifcation. In: Thirty-Second AAAI Conference on Artifcial Intelligence
- <span id="page-27-18"></span>25. He K, Zhang X, Ren S, Sun J (2016) Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778
- <span id="page-27-19"></span>26. Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classifcation with deep convolutional neural networks. Adv Neural Inf Process Syst 25:1097–1105
- <span id="page-27-20"></span>27. Kim C, Li F, Ciptadi A, Rehg JM (2015) Multiple hypothesis tracking revisited. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 4696–4704
- <span id="page-27-21"></span>28. Girshick R (2015) Fast r-cnn. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 1440–1448
- <span id="page-28-0"></span>29. Long J, Shelhamer E, Darrell T (2015) Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3431–3440
- <span id="page-28-1"></span>30. Bolya D, Zhou C, Xiao F, Lee YJ (2019) Yolact: Real-time instance segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 9157–9166
- <span id="page-28-2"></span>31. Zheng C, Cham T-J, Cai J (2019) Pluralistic image completion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1438–1447
- <span id="page-28-3"></span>32. Ji R, Chen F, Cao L, Gao Y (2018) Cross-modality microblog sentiment prediction via bi-layer multimodal hypergraph learning. IEEE Trans Multim 21(4):1062–1075
- <span id="page-28-5"></span>33. Feng Y, You H, Zhang Z, Ji R, Gao Y (2019) Hypergraph neural networks. In: Proceedings of the AAAI Conference on Artifcial Intelligence, vol. 33, pp. 3558–3565
- <span id="page-28-4"></span>34. Bai S, Zhang F, Torr PH (2021) Hypergraph convolution and hypergraph attention. Pattern Recognit 110:107637
- <span id="page-28-6"></span>35. Pohlen T, Hermans A, Mathias M, Leibe B (2017) Full-resolution residual networks for semantic segmentation in street scenes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4151–4160
- <span id="page-28-9"></span>36. Krishna R, Zhu Y, Groth O, Johnson J, Hata K, Kravitz J, Chen S, Kalantidis Y, Li L-J, Shamma DA, Bernstein M, Fei-Fei L (2016) Visual genome: Connecting language and vision using crowdsourced dense image annotations.<https://arxiv.org/abs/1602.07332>
- <span id="page-28-7"></span>37. Duan K, Bai S, Xie L, Qi H, Huang Q, Tian Q (2019) Centernet: Keypoint triplets for object detection. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 6569–6578
- <span id="page-28-8"></span>38. Mikolov T, Chen K, Corrado G, Dean J (2013) Efficient estimation of word representations in vector space. <http://arxiv.org/abs/1301.3781>
- <span id="page-28-10"></span>39. Wu Y, He K (2018) Group normalization. In: Proceedings of the European Conference on Computer Vision (ECCV), pp. 3–19
- <span id="page-28-11"></span>40. Kingma DP, Ba J (2014) Adam: A method for stochastic optimization. [http://arxiv.org/abs/1412.](http://arxiv.org/abs/1412.6980) [6980](http://arxiv.org/abs/1412.6980)
- <span id="page-28-12"></span>41. Cordts M, Omran M, Ramos S, Rehfeld T, Enzweiler M, Benenson R, Franke U, Roth S, Schiele B (2016) The cityscapes dataset for semantic urban scene understanding. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3213–3223
- <span id="page-28-13"></span>42. Hariharan B, Arbeláez P, Girshick R, Malik J (2014) Simultaneous detection and segmentation. In: European Conference on Computer Vision, pp. 297–312. Springer
- <span id="page-28-14"></span>43. Yang T-J, Collins MD, Zhu Y, Hwang J-J, Liu T, Zhang X, Sze V, Papandreou G, Chen L-C (2019) Deeperlab: Single-shot image parser.<http://arxiv.org/abs/1902.05093>

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional afliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.

## **Authors and Afliations**

# **Chang‑Eun Lee<sup>1</sup> · Jaeuk Baek1 · Jeany Son<sup>2</sup> · Young‑Guk Ha3**

Chang-Eun Lee celee@etri.re.kr

Jaeuk Baek jubaek@etri.re.kr

Jeany Son jeany@gist.ac.kr

<sup>1</sup> Electronics and Telecommunications Research Institute (ETRI), Daejeon, Republic of Korea

- <sup>2</sup> Gwangju Institute of Science and Technology (GIST), Gwangju, Republic of Korea
- <sup>3</sup> Konkuk University, Seoul, Republic of Korea