

# Personalized Saliency in Task-Oriented Semantic Communications: Image Transmission and Performance Analysis

Jiawen Kang, Hongyang Du, Zonghang Li, Zehui Xiong, Shiyao Ma, Dusit Niyato, *Fellow, IEEE*, Yuan Li

**Abstract**—Semantic communication, as a promising technology, has emerged to break through the Shannon limit, which is envisioned as the key enabler and fundamental paradigm for future 6G networks and applications, e.g., smart healthcare. In this paper, we focus on UAV image-sensing-driven task-oriented semantic communications scenarios. The majority of existing work has focused on designing advanced algorithms for high-performance semantic communication. However, the challenges, such as energy-hungry and efficiency-limited image retrieval manner, and semantic encoding without considering user personality, have not been explored yet. These challenges have hindered the widespread adoption of semantic communication. To address the above challenges, at the semantic level, we first design an energy-efficient task-oriented semantic communication framework with a triple-based scene graph for image information. We then design a new personalized semantic encoder based on user interests to meet the requirements of personalized saliency. Moreover, at the communication level, we study the effects of dynamic wireless fading channel on semantic transmission mathematically and thus design an optimal multi-user resource allocation scheme by using game theory. Numerical results based on real-world datasets clearly indicate that the proposed framework and schemes significantly enhance the personalization and anti-interference performance of semantic communication, and are also efficient to improve the communication quality of semantic communication services.

**Index Terms**—Semantic communication, personalized saliency, resource allocation, unmanned aerial vehicle.

## I. INTRODUCTION

### A. Background and Motivations

The fast-growing 6G communication technology is enabling the transition from serving people and things to supporting the “Internet of Everything” [1]–[3]. Specifically, 6G communication technology serves intelligent production and life through intelligent interconnection and collaborative symbiosis of human-machine-object, and actively promotes the construction of an inclusive and intelligent human society [4]–[6]. In the 6G era, however, the traditional point-to-point information

transmission communication system, that relies on the resource optimization of physical-layer dimension and the stable transmission protocol at the network layer, cannot meet the increasing requirements of complex, diverse, and intelligent information transmission needs, e.g., supporting virtual reality, holographic projection and Metaverse applications [7], [8]. Therefore, it is essential to design a new communication paradigm for efficient information transmission thus meeting the demands of future communication.

Fortunately, semantic communication [9], as a new architecture that can integrate user needs and information semantic features into the communication process, is expected to become a new communication paradigm for the Internet of Everything in the future [10]–[12]. Different from the traditional communication architectures, the semantic communication system aims to fundamentally solve the problems of cross-system, cross-protocol, cross-network, and cross-human-machine information transmission redundancy in traditional information-transmission based communication protocols. The ultimate goal of the semantic communication system is to efficiently transmit content-aware and semantic-related information in a task-oriented manner, and make the grand vision of “Internet of Everything” come true. In this context, to better serve the two core technical standards of semantics and validity in 6G networks, the Task-Oriented Semantic Communication Systems (TOSC) was designed with two parts: semantic reconstruction and goal execution [13]

For semantic reconstruction, a semantic feature extractor is applied to extract the semantic features behind the data to be transmitted and reconstruct the semantic information at the receiver. For example, Zhu *et al.* in [14] proposed an adaptive transformer to encode the semantic information and decode it at the receiver. For task-oriented applications, semantic communication aims to extract semantic information related to the decision goal of the receiver. Prior work generally focuses on image recognition scenarios and develops image classification-oriented semantic communication for improving recognition accuracy rather than image semantic reconstruction [15], [16]. Since Unmanned Aerial Vehicle (UAV) sensing has been widely applied in various industries with 6G, in this paper, we focus on investigating UAV sensing-driven task-oriented semantic communication systems.

Particularly, UAV-sensing-driven task-oriented semantic communication aims to provide users with cross-regional intelligent services with the help of UAV’s image retrieval, image recognition, image transmission, and image coding features.

Jiawen Kang is with Automation of School, Guangdong University of Technology, China. (Email: kavinkang@gdut.edu.cn).

Hongyang Du and Dusit Niyato are with School of Computer Science and Engineering, Nanyang Technological University, Singapore.

Zonghang Li is with School of Information and Communication Engineering, University of Electronic Science and Technology of China, Chengdu, China. (Email: zhli@std.uestc.edu.cn).

Zehui Xiong is with Pillar of Information Systems Technology and Design, Singapore University of Technology and Design, Singapore.

Shiyao Ma is with College of Information and Communication Engineering, Dalian Minzu University.

Yuan Li is with Shangdong University, China.

UAV-sensing-driven Task-oriented Semantic Communication (UTSC) generally serves a multi-demand, complex cross-modal intelligent task with multiple users. The UTSC is particularly suitable for emerging intelligent scenarios in daily life [6], [9] and industry [17] (e.g., smart agriculture [18]). The service mode of UTSC is realized by collecting task demands from users and using UAV sensing equipment to collect image data, and returning demand feedback to users through cloud servers in the form of semantic information. For example, in smart agriculture, different farmers (i.e., users) require UAVs to accomplish different tasks (e.g., monitoring and shooting). In this case, it is challenging for UAVs to accomplish the task goals of all users simultaneously. Furthermore, due to resource limitations and power supply problems, UAVs cannot hover in the air for a long time to complete the personalized user demands in turn. Therefore, there exists following unique challenges when designing USTC systems:

- **Energy-hungry and efficiency-limited image retrieval:** Traditional communication manners generally send all the captured images to the users, while many of the images are not interesting/needed by the users. Therefore, such inefficient communication manners may cause large waste of communication resources of UAVs and consume unnecessary UAV energy. A straightforward solution is to use keyword subscriptions, namely, UAVs push images of interest to users by extracting scene graphs of sensing images that match the subscription words provided by users. However, this solution cannot execute well in the traditional communication manners under the scenarios of limited UAV energy or poor wireless channel. It is because that image packet drop and re-transmission bring large delay. Thus the inefficiency of wireless transmission makes it difficult to realize real-time push according to keyword subscriptions. To this end, it is urgently necessary to design an efficient and energy-efficient semantic-based real-time subscription method to improve subscription accuracy and wireless resource utilization.
- **Semantic encoding without personality:** Existing work ignores the personalized needs of semantic communication for encoding retrieval results and does not set user preferences for the importance of semantically encoded values. In the face of large-scale semantic communication, due to the limited bandwidth resources between UAVs and users, semantic coding is prone to signal fading during transmission, resulting in the drop of important coding values (i.e., important information for users). Therefore, we need to design a personalized semantic coding value weight setting scheme, and design different resource allocation schemes for users with different interests to ensure that important information about user preferences can be efficiently conveyed.
- **Without insights between wireless fading channel and semantic communication:** The semantic triplet drop probability is defined as the probability that the number of error bits in the triplet exceeds the error correction capability. However, existing works rarely consider the drop of encoded information caused by the physical

environment during information transmission. A common approach is to use neural networks to model the wireless channel, which cannot help the system design with mathematical analysis [19]. The problems about the multi-path effect, shadow effect, and co-channel interference for wireless transmission environment affect the communication quality and increase the semantic triplets drop problem during transmission. Thereby, it is essential to study the joint optimization problem of the wireless fading channel and semantic triplet drop probability, and also design a resource optimization scheme.

## B. Solutions and Contributions

To address the above challenges, in this paper, we propose an energy-efficient task-oriented semantic communication framework for 6G-enabled UAV image sensing scenarios. More specifically, in this framework, image information is modeled as a triple-based scene graph to provide users with images that meet their preference requirements in an efficient retrieval manner. On this basis, we further execute the triplets by weight-encoding and use a personalized attention-based mechanism to implement differential weight encoding of triplets for important information according to user preferences. Moreover, for the UAV power allocation issues, we further consider the dynamic wireless fading effects on the semantic information transmission, and thus mathematically analyze the triplet drop probability. Based on the theoretical analysis results, we formulate a game theory model and design a multi-user resource allocation scheme to achieve efficient resource utilization and maximize the resource utility of UAVs.

The key contributions are summarized as follows:

- Unlike traditional UAV-sensing communications that require all images to be transmitted, we propose an energy-efficient semantic communication-based framework that enables UAV only to transfer the selected interested images of the users, which is achieved by matching the user's query text with the semantic information of all images.
- We design a novel personalized semantic encoder with the help of the user's subjective interest. After obtaining the semantic information, i.e., triples, from the image, the triplets of more interest to the user are given higher weights, thus ensuring the correct reception.
- We analyze the performance of semantic triplets transmission mathematically. We derive the exact expression for the semantic triplet drop probability by considering the generalized fading channel model. From the derived expressions, we obtain insights into the wireless channel environment on the impact of semantic communication.
- Considering the resource limitation of UAVs and the requirements to send semantic information to multiple users, we propose a multi-user resource allocation scheme based on game theory, whose utility function of the retrieval task is used as the optimization objective for better resource utilization.

The remainder of this paper is organized as follows. We first summarize the related work about Semantic Communication

TABLE I  
SUMMARY OF MAIN SYMBOLS.

Symbol	Explanation
$m_f$	The fading parameter for the $k_{\text{th}}$ user in Fisher-Snedecor $\mathcal{F}$ fading model.
$m_s$	The shadowing parameter for the $k_{\text{th}}$ user in Fisher-Snedecor $\mathcal{F}$ fading model.
$\bar{z}$	The average value of $\mathcal{F}$ random variables (RVs), i.e., $z$ , for the $k_{\text{th}}$ user in Fisher-Snedecor $\mathcal{F}$ fading model.
$\mathbf{a}^T$	Transpose of vector $\mathbf{a}$ .
$\Gamma(\cdot)$	Gamma function [20, eq. (8.310.1)].
$\Gamma(\cdot, \cdot)$	Upper incomplete Gamma function [20, eq. (8.350.2)].
$B(\cdot, \cdot)$	Beta function [20, eq. (8.384.1)].
$F(\cdot, \cdot; \cdot; \cdot)$	Gauss hypergeometric function [20, eq. (9.111)], which is also known as ${}_2F_1(\cdot, \cdot; \cdot; \cdot)$ .
$G_{p,q}^{m,n}(\cdot)$	Meijer's $G$ -function [20, eq. (9.301)].
$H_{\cdot}(\cdot, \cdot)$	Multivariate Fox's $H$ -function [21, eq. (A-1)].
$G_{\cdot, \cdot}(\cdot, \cdot)$	Bivariate Meijer's $G$ -function [22, eq. (1)].
$N_T$	The number of antennas in the UAV
$x$	The input image.
$K$	The number of users.
$k$	The identity of user.
$q^k$	The personal query text of user $k$ .
$\mathcal{T}$	The Triplet Detection (TD) function.
$\mathcal{P}$	The Personalized Saliency Prediction (PSP) function.
$\mathcal{C}$	The crop function, used to crop a sub-image from $x$ according to the box coordinates.
$\tau$	The set of (subject-relation-object) triplets.
$\tau_{\text{recv}}^k$	The triplets received by user $k$ .
$H_{\text{sub}}, H_{\text{obj}}$	The attention heatmaps of subject and object entities.
$B_{\text{sub}}, B_{\text{obj}}$	The box coordinates of subject and object entities.
$p^k$	The personalized triplet priority of user $k$ .
$S^k$	The personalized saliency heatmap of user $k$ .
$\alpha$	The coefficient for fusing attention heatmap and saliency heatmap.
$F_{\text{sub}}^k, F_{\text{obj}}^k$	The fused attention heatmaps of subject and object entities of user $k$ .
$\tilde{F}_{\text{sub}}^k, \tilde{F}_{\text{obj}}^k$	The sub-heatmaps of subject and object entities of user $k$ , cropped from $F_{\text{sub}}^k, F_{\text{obj}}^k$ .
$s$	The match score between the triplets received and the personal query.

in Sec. II. We then introduce our system model in Sec. III. Then, we present the proposed semantic triplets transmission method from the communication level in Sec. IV. Next, we elaborate on the proposed semantic communication framework in Sec. V. Subsequently, we conduct a series of case studies and analyze the simulation results in Sec. VI. Finally, we conclude this paper in Sec. VII. We summarize the mathematical symbols and explanations in Table I.

## II. RELATED WORK

### A. Semantic Communication

Semantic communication is a new communication paradigm, and it can transmit more information when the external environment is the same [23]. The sender sends the semantic information extracted from the original information (such as images, text, and video) to the receiver in a semantic communication system and the receiver recovers the original information from the semantic information. By transmitting the semantic information, the semantic communication system can eliminate the unnecessary information from the original information to reduce communication overhead. The research on semantic communication is roughly divided into four parts: how to optimize semantic encoding, how to complete the goal-oriented communications, how to protect semantic information privacy, and how to analyze the systems' performance.

**Semantic Encoding Optimization:** Semantic encoding is crucial for the semantic communication system because it determines the efficiency and effectiveness of information transmission. Xie et al. [19] proposed a DeepSC framework that is based on a Transformer for text transmission tasks, DeepSC can extract the semantic information from the original text under the interference of noise. Xie et al. [17] further considered the more realistic situation that designing a lightweight deep learning semantic communication system makes the model easier to deploy on IoT devices. The above work only focused on text-based semantic communication. Besides, images-based semantic communication has been proposed. Bourtsoulatze [24] proposed a semantic communication system that utilizes joint source and channel coding techniques to reduce communication overhead for image transmission tasks. Although the above work [17], [19], [24] proposed the semantic communication system to reduce the communication overhead, they only considered the transmission on single-domain, such as text-to-text or image-to-image. It is unrealistic in the real world. Therefore, we consider the transmission of multi-domain and multimodal in this paper. We design the framework in which the sender extracts the triplets with semantic information from the original images, and then the sender sends the triplets to the receiver. It can further reduce the communication overhead and be more practical.

**Task-oriented Communications:** The task-oriented communication is also important for the semantic communication system because it transmits different semantic information according to different needs. That is to say, task-oriented communication systems do not need to transmit data directly under corresponding circumstances. To this end, Farshbafan et al. [25] proposed a task-oriented semantic communication framework, which can transmit different semantic information in the different goals of the system. However, reference [25] only considered a single sender and receiver at the same time. Xie et al. [13] proposed a task-oriented multi-user semantic communication system, the framework utilized different encoders and decoders to solve the multi-task and multi-user problems. Although the above work [13], [25] proposed the task-oriented communication system to complete

the corresponding goal, they ignored the individual differences in different users, i.e., other receivers will have different needs in the real world. Therefore, we consider the personalized saliency in our framework to adapt to different users' goals.

**Secure Semantic Encoding:** Since semantic information can reflect the real data distribution of users to a certain extent and is also vulnerable to privacy leakage in communication, we need to protect the semantic information transmitted by users. For example, Chen et al. [26] proposed a federated learning framework in semantic communication systems, which can protect the privacy of the system. Yang et al. [27] proposed a federated learning-based semantic communication framework to cope with the computing and communication overhead under protecting the privacy of the system. The above work illustrates that federated learning can be applied to wireless communication systems to protect the devices' data privacy. Therefore, Tong et al. [28] designed a wav2vec-based autoencoder federated semantic communication framework, and it can significantly reduce transmission error and heavy communication overhead. Due to the low efficiency of the current semantic communication system, we focus on how to extract semantic information effectively and effectively adapt to different goals.

**Semantic Communications Performance Analysis:** In a wireless semantic communication system, the transmission performance of the system is negatively affected by multipath fading, that is, interference between different signals. Therefore, it is crucial to accurately model fading channels to better understand the performance impact of fading channels on semantic communication systems. Previous work focused on the performance analysis of wireless communication systems. For example, Yoo et al. [29] proposed the  $\mathcal{F}$  distribution as fading model to analyze the performance of the semantic communication systems. Based on this, reference [30] further explored a comprehensive performance analysis of the  $\mathcal{F}$  composite fading channels in conventional wireless communication systems. Due to the rise of semantic communications, some researchers currently turned attention to studying the impact of fading channels on the performance of wireless semantic communication systems. Xie et al. [31] and Weng et al. [32] explored the performance of their framework over Rayleigh and Rician channels in semantic communication systems. However, reference [31], [32] only explored the communication performance in Rayleigh and Rician fading channels that are the traditional multipath fading models. Due to the Fisher-Snedecor  $\mathcal{F}$  fading channel is a commonly used fading channel model in wireless communication systems, and it can cover the situation of classical fading channel analysis through changing parameters. Therefore, in this paper, we analyze the performance of our model by considering Fisher-Snedecor  $\mathcal{F}$  channel model in wireless semantic communication and obtain insights into the wireless channel environment on the impact of semantic communication.

### B. Personalized Saliency

Personalized saliency means that different observers have different regions of interest in the same image. In recent

years, personalized saliency has received a lot of attention in the computer vision community. Xu et al. [33] proposed a CNN-based and personal information framework to predict a personalized saliency map. However, personalized information and images always change in the real world. Dodge et al. [34] proposed a framework that combines the global scene information from all categories and the extracted local information to predict the saliency map. However, many categories are unknowable. Mahdi et al. [35] utilized three CNN-based models to obtain the bottom-up and top-down deep features to predict a personalized saliency map. However, due to the storage capacity and the computing resources, it may be difficult to utilize this framework for extracting complex features when the number of users increases. Berkovsky et al. [36] proposed a framework for predicting personalized saliency based on eye tracking data and the framework needs to capture physiological responses, such as brain signals. However, there exist massive users in semantic communication systems and it is difficult to capture signals between multiple users, so the framework is difficult to adapt to multi-user situations. Moroto et al. [37] proposed the framework to extract personalized saliency maps (PSMs) through Gaussian process regression. However, considering the limitation of storage capacity, the receivers, i.e., UAVs, are difficult to store the PSMs for all users. Therefore, we convert PSMs into triples to further reduce memory overhead in this paper. The above work may be unsuitable for semantic communication systems, especially for UAV scenarios. Therefore, we consider the situation of a multi-user semantic communication system and propose personalized saliency-based semantic communication.

In summary, we propose a personalized saliency-based task-oriented semantic communication system to cope with the above problems in this paper. Firstly, we predict the saliency heatmap of the user through the customized information and the image captured by the UAV. Meanwhile, UAV executes triplet detection [38] to generate an attention heatmap from the image captured by the UAV. Secondly, UAV executes the attention fusion step to obtain the fused attention for each user and obtains the personalized triplet from the fused attention. The above steps complete the purpose of personalized saliency. Thirdly, UAV allocates multi-user power through triplet priority estimation and transmits the triplets by power allocated. Finally, the user obtains the match score between the received triplets and personal query and decides whether to download the image based on the matching score. The last steps complete the purpose of goal-oriented semantic communication.

## III. SYSTEM MODEL

In this section, we first describe our proposed task-oriented SemCom system and then describe the metrics used to evaluate the effectiveness of the proposed design.

### A. Task-oriented SemCom Design

When users hire a UAV for image acquisition, in many cases, not all the images taken by the UAV are what the user needs. Therefore, a retrieval task needs to be performed on all the images, i.e., the users input the text of the scene they

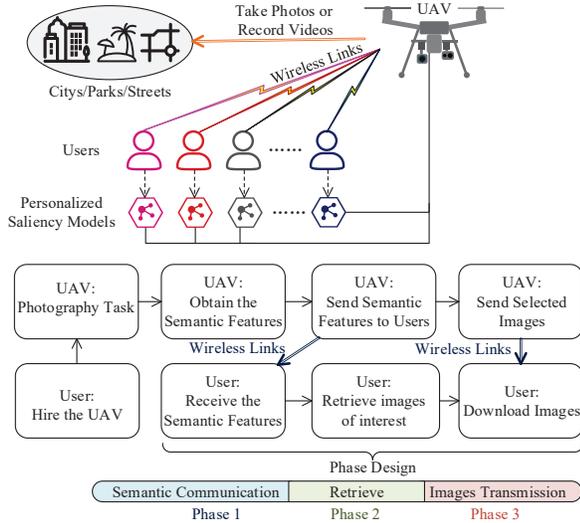


Fig. 1. An illustration of the proposed task-oriented SemCom system model.

want and get the corresponding images. However, considering that the energy of the UAV is limited and that one UAV may need to serve multiple users, performing all the users' retrieval tasks in the UAV will affect the quality of images and the UAV's endurance. Moreover, sending all the images to users and retrieving images on the user side also involve unnecessary energy consumption, i.e., users need to download some images that are not of interest.

The development of semantic communication has given us a solution to solve the above problems. As shown in Fig. 1, instead of transmitting all the images back to users after the photography task, the UAV transmits the semantic features (text format) of the images to users. The transmission of semantic features of images in text format requires few channel resources and is convenient for users to store. The users determine a query text according to their interests, match the images in the received semantic features, and then download the original image from the UAV. In the rest of this paper, we focus on the extraction and transmission of semantic information.

## B. Optimization Problems Analysis

The UAV extracts the triplets that represent the corresponding image information from the captured images. If  $n_i$  triplets can be extracted from the  $i_{\text{th}}$  image, we have  $\tau_i = \{\tau_{i,1}, \dots, \tau_{i,n_i}\}$ , where  $\tau_i$  denotes the set of triplets from the  $i_{\text{th}}$  image. In Phase 1, the UAV transmits the triplets to the users. Note that users are not necessarily located close to each other, and the different interests of users will make them have different query texts. Therefore, we consider the TDMA scheme in Phase 1. In each time slot, i.e.,  $T_1$ , the UAV uses all the antennas to serve one user and designs the beamforming vector accordingly. Let us consider an energy-constrained scenario, which means that

$$\sum_{j \in \mathcal{K}} P_j T_1 < W_A, \quad (1)$$

where  $\mathcal{K}$  is the users selected in this round,  $P_j$  is the transmit power for the  $j_{\text{th}}$  user and  $W_A$  is the total energy. Thus, we need to solve two optimization problems:

- **P1-power allocation among users:** One user per time slot is served. Because the total energy of the UAV is limited, we need to determine how many resources each user can occupy, in the form of transmit power.
- **P2-power allocation among one user's triplets:** After determining the available power for each user, we also need to determine the power that should be allocated to each triplet.

Let  $s_k$  denote the match score between the triplets received and the personal query of the  $k_{\text{th}}$  user,

$$s_k = \frac{N_{\text{in},k}}{N_{\text{rec}}}, \quad (2)$$

which represents the proportion of the number of interested images of the  $k_{\text{th}}$  user obtained by matching, i.e.,  $N_{\text{in},k}$ , to the total number of images captured by UAV, i.e.,  $N_{\text{rec}}$ .

Let  $\tilde{s}_k$  denote the optimal match score,

$$\tilde{s}_k = \frac{N_{\text{truth},k}}{N_{\text{rec}}}, \quad (3)$$

which represents the proportion of the number of interested images of the  $k_{\text{th}}$  user in truth (without considering the impact of triplets drop that may be caused by wireless transmission), i.e.,  $N_{\text{truth}}$ , to  $N_{\text{rec}}$ . Then, we can use  $\frac{s_k}{\tilde{s}_k}$  to represent the effectiveness of semantic communication.

If P1 can be solved, we derive the power resources available to each user. Then, P2 for each user can be solved with the help of the personalized weights for the triplets. The solution for P2 is discussed in Section V. Now we focus on P1. We consider the power allocation problem among users as cooperative bargaining. Thus, with the help of NBS, we need to maximize  $\prod_{k=1}^K \frac{s_k}{\tilde{s}_k}$ . Because  $\tilde{s}_k$  represents the ground-of-truth which is only decided by the query text, the optimization problem can be transformed into

$$\max_{P_k(k=1, \dots, K)} \prod_{k=1}^K s_k. \quad (4)$$

For a given query text of the  $k_{\text{th}}$  user, the  $s_k$  is mainly affected by the number of received triplets. Due to the random nature of the wireless communication environment, the transmitted triplets may not be decoded correctly because of excessive error codes. Therefore, the higher transmit power should be allocated to the more important triplets, i.e., the triplets that are of more interest to the user, to guarantee error-free transmission.

## C. SINR Analysis

We consider a set of  $\mathcal{K} = 1, \dots, K$  user, where each user has their own preference. One UAV performs the photography task and needs to transmit the semantic features of the image to  $K$  users. To obtain considerable array gains and improve the channel quality, we consider the UAV is equipped with  $N_T$  antennas. The UAV is hovering above the ground  $K$

users. The horizontal coordinate of the  $k_{\text{th}}$  ground device is assumed to be  $u_k = (x_k, y_k, 0)$ ,  $k = (1, \dots, K)$ , while the UAV is hovering at a fixed altitude  $z_u$  with the coordinate  $u_u = (x_u, y_u, z_u)$ . Thus, the distance between the  $k_{\text{th}}$  user and the UAV can be expressed as  $D_{uk} = \sqrt{\|u_k - u_u\|^2}$ . Let  $\alpha_k$  denote the path loss exponents of the UAV- $k_{\text{th}}$  user link.

We denote the channel vector from the UAV to the  $k_{\text{th}}$  user as  $\mathbf{h}_k \in \mathbb{C}^{1 \times K}$ . By adopting the linear beamforming, the data symbol  $t_k$  intended for user  $k$  is multiplied with the beamformer  $\mathbf{w}_k \in \mathbb{C}^{K \times 1}$ .  $N_{Ik}$  paths of interferes are assumed to be present at the  $k_{\text{th}}$  user, where  $I$  means the abbreviation of interference, which is used to distinguish symbols. Each of the  $j_{\text{th}}$  user interfering signals has an average transmit power  $P_{Ik}$ .  $t_{I,k,j}$  is the  $j_{\text{th}}$  interfering symbol. Accordingly, because TDMA is used, the received signal at the  $k_{\text{th}}$  user is given by<sup>1</sup>

$$r_k = \sqrt{P_k D_{uk}^{-\alpha_k}} \mathbf{h}_k \mathbf{w}_k x_k + P_{Ik} \sum_{j=1}^{N_{Ik}} \sqrt{h_{I,k,j}} t_{I,k,j} + n_k, \quad (5)$$

where  $n_k \in \mathcal{CN}(0, \sigma^2)$  is the noise,  $P_k$  is the transmit power for the  $k_{\text{th}}$  user. The SINR at the  $k_{\text{th}}$  user can be expressed as

$$\gamma_k = \frac{P_k D_{uk}^{-\alpha_k} \|\mathbf{h}_k \mathbf{w}_k\|^2}{\sigma^2 + P_{Ik} \sum_{j=1}^{N_{Ik}} h_{I,k,j}^2}. \quad (6)$$

With the help of maximum ratio transmission [39], the optimal beamforming vector can be expressed as  $\mathbf{w}_k = \frac{\mathbf{h}_k^T}{\|\mathbf{h}_k\|}$ . Thus, we have  $\|\mathbf{h}_k \mathbf{w}_k\|^2 = \sum_{j=1}^{N_T} h_{k,j}^2$ .

### D. Channel Model

The small scale fading of UAV- $k_{\text{th}}$  user link is modeled as the Fisher-Snedecor  $\mathcal{F}$  fading distribution. The Fisher-Snedecor  $\mathcal{F}$  composite fading model assumes that small-scale variations follow the Nakagami- $m$  distribution and shadowing follows the inverse Nakagami- $m$  distribution. Channel measurements at 5.8 GHz have demonstrated that the Fisher-Snedecor  $\mathcal{F}$  fading model fits experimental results better than the KG fading model both in line-of-sight and non-LOS scenarios [29].

Thus,  $\|\mathbf{h}_k \mathbf{w}_k\|^2$  follows the distribution of sum of  $N_T$  Fisher-Snedecor  $\mathcal{F}$  RVs [40]. However, the PDF and CDF of  $\|\mathbf{h}_k \mathbf{w}_k\|^2$  is in terms of Multivariate Fox's  $H$ -function [21, eq. (A-1)], which is hard to provide insights. Considering that the Fisher-Snedecor  $\mathcal{F}$  RV is defined as the ratio of two Gamma RVs and the sum of Gamma RVs still follows the Gamma distribution, we can use the single Fisher-Snedecor  $\mathcal{F}$  distribution to approximate the distribution of the sum of Fisher-Snedecor  $\mathcal{F}$  RVs [40].

Let  $Z \triangleq \|\mathbf{h}_k \mathbf{w}_k\|^2 \sim \mathcal{F}(m_{fk}, m_{sk}, \bar{z}_k)$ , the PDF and CDF of  $Z$  are given as [41, eq. (6)] and [41, eq. (12)], respectively.

<sup>1</sup>Note that the large scale fading of the interference signal is considered in the mean value of  $h_{I,k,j}$ .

We consider the interference signals follow the Rayleigh distribution, i.e.,  $h_{I,k,j} \sim \text{Rayleigh}(\eta_k)$ . Let  $Y \triangleq \sum_{j=1}^{N_{Ik}} h_{I,k,j}^2$ . Because the sum of  $N_{Ik}$  i.i.d. Rayleigh-fading signals have a Nakagami- $m$  distributed signal amplitude with  $m = N_{Ik}$ , the PDF and CDF expressions of  $\|\mathbf{h}_k\|^2$  can be written as [42].

$$f_Y(y) = \frac{y^{N_{Ik}-1}}{\eta_k^{N_{Ik}} \Gamma(N_{Ik})} \exp\left(-\frac{y}{\eta_k}\right), \quad (7)$$

and

$$F_Y(y) = \frac{\Gamma\left(N_{Ik}, \frac{y}{\eta_k}\right)}{\Gamma(N_{Ik})}, \quad (8)$$

where  $\eta_k = \mathbb{E}[h_k^2]$ , and  $\mathbb{E}[\cdot]$  denotes expectation. We then derive the PDF and CDF of the SINR,  $\gamma_k = \frac{P_k D_{uk}^{-\alpha_k} Z}{\sigma^2 + P_{Ik} Y}$ .

**Theorem 1.** *The PDF and CDF of the SINR,  $\gamma_k = \frac{P_k D_{uk}^{-\alpha_k} Z}{\sigma^2 + P_{Ik} Y}$ , can be derived as (9) and (10), respectively, shown at the bottom of the next page, where  $\Lambda_k \triangleq \frac{P_k D_{uk}^{-\alpha_k} (m_{sk}-1)\bar{z}_k}{m_{fk}}$ ,  $\Theta_1 = (1 - N_{Ik} : \{-1, -1, 0\})$ , and  $\Theta_2 = (1 - m_{fk} - m_{sk} : \{0, -1, -1\})$ .*

*Proof:* Please refer to Appendix A. ■

### E. Approximation Analysis

Although the previously derived PDF and CDF are obtained in closed-form, it is hard to bring valuable insights if we derive performance expressions with the help of (9) and (10). Therefore, in the following, we derive the accurate approximate CDF by presenting an approximate solution to a complex mathematical integral equation. Moreover, we analyze the approximate CDF in the high-SNDR regime and verify our derived results by numerical analysis.

1) *Accurate Approximation:* Let us consider the integration as follows:

$$I_A = \int_a^\infty x^b (x-a)^c \exp\left(\frac{x-a}{d}\right) F(\alpha, \beta, \varepsilon; \rho x) dx. \quad (11)$$

Note that  $I_A$  has not been studied in any mathematical integral theory book or website such as [20], [43]. It is difficult, if not impossible, to obtain the closed solution of  $I_A$ . Here we derive the accurate approximate solution of  $I_A$  in the Lemma 1.

**Lemma 1.** *An accurate approximation of  $I_A$  when  $\rho$  is small can be expressed as*

$$I_A = \exp\left(\frac{a}{d}\right) \frac{\Gamma(\varepsilon) d^{b+c+1}}{\Gamma(\alpha) \Gamma(\beta)} G_{3,2}^{1,3} \left( d\rho \left| \begin{matrix} 1-\beta, -b-c, 1-\alpha \\ 1-\varepsilon \end{matrix} \right. \right). \quad (12)$$

*Proof:* Please refer to Appendix B. ■

With the help of Lemma 1, we can re-derive the CDF of  $\gamma_k$  as

**Theorem 2.** *An accurate CDF of  $\gamma_k$  can be obtained as*

$$F_{\gamma_k}(\gamma) = \frac{1}{\Gamma(N_{Ik}) \Gamma(m_{sk}) \Gamma(m_{fk})} \exp\left(\frac{\sigma^2}{P_{Ik} \eta_k}\right) \times G_{3,2}^{1,3} \left( \frac{\gamma P_{Ik} \eta_k}{\Lambda_k} \left| \begin{matrix} 1 - m_{sk}, 1 - N_{Ik}, 1 \\ m_{fk}, 0 \end{matrix} \right. \right). \quad (13)$$

2) *High-SNDR approximation*: In the following, we analyze asymptotic CDF in the high-SNDR regime.

**Theorem 3.** *The CDF of  $\gamma_k$  can be approximated in the high transmit power regime as*

$$F_{\gamma_k}(\gamma) = \frac{\Gamma(m_{fk} + m_{sk}) \Gamma(N_{Ik} + m_{fk})}{\Gamma(N_{Ik}) \Gamma(m_{sk}) \Gamma(m_{fk} + 1)} \times \exp\left(\frac{\sigma^2}{P_{Ik}\eta_k}\right) \left(\frac{\gamma P_{Ik}\eta_k}{\Lambda_k}\right)^{m_{fk}}. \quad (14)$$

3) *Verification and Insights*: We use the OP to verify the derived CDF expression, (10), and two approximate expressions, (12) and (14). The outage probability (OP) is defined as the probability that the SINR falls below a given outage threshold, i.e.,  $OP_k = \mathbb{P}(\gamma_k < \gamma_{th}) = F_{\gamma_k}(\gamma_{th})$ . As shown in Fig. 2, we study the OP versus the transmit power, with  $D_{uk} = 1.5$  m,  $\alpha_k = 2$ ,  $m_{sk} = 5$ ,  $m_{fk} = 2.6$ ,  $\bar{z}_k = -1$  dB,  $\sigma^2 = 1$  W,  $\eta_k = 0.4$ ,  $P_{Ik} = 5$  W,  $N_{Ik} = 3$ , and different values of  $m_{fk}$ . We can observe that the accurate approximate expression (12) matches almost exactly with the closed-form analytic expression (10). In the high-SINR regime, e.g., when  $P_k$  is larger than 25 dBW, the values obtained from the high-SINR approximate expression (14) are close to that obtained from analytic expression. Furthermore, from (14), we can observe that the multi-path fading parameter  $m_{fk}$ , instead of the shadowing parameter  $m_{sk}$ , determines the slope of the OP with decreasing transmit power. In Fig. 2, it is shown that the larger the  $m_{fk}$  is, the more rapidly the OP decreases with the increase of transmit power.

#### IV. COMMUNICATION LEVEL: TRIPLET DROP PROBABILITY

We encode the semantic information of the image into triplets. Due to the instability of wireless transmission, the receiver cannot guarantee the perfect receiving of the semantic triplets transmitted by the UAV. Therefore, we analyze the impact of the wireless environment on the semantic triplets transmission. Suppose a triplet is encoded with bit length  $D_T$ , and that the use of bit error correction codes allows for at most  $D_E$  error bits.

##### A. Bit Error Probability

Under a variety of modulation formats, the BEP,  $E_k$ , can be expressed as [44, eq. (13)]

$$E_k = \int_0^\infty \frac{\Gamma(\lambda_2, \lambda_1 \gamma)}{2\Gamma(\lambda_2)} f_{\gamma_k}(\gamma) d\gamma, \quad (15)$$

where  $\Gamma(\lambda_2, \lambda_1 \gamma)/2\Gamma(\lambda_2)$  is the conditional bit error probability,  $\lambda_1$  and  $\lambda_2$  are modulation-specific parameters which have different values under different modulation and detection schemes [44].

**Theorem 4.** *The BEP of the  $k_{th}$  user can be derived as*

$$E_k = \frac{\Gamma^{-1}(m_{sk}) \Gamma^{-1}(m_{fk})}{2\Gamma(\lambda_2) \Gamma(N_{Ik})} \exp\left(\frac{\sigma^2}{P_{Ik}\eta_k}\right) \times G_{4,2}^{1,4} \left( \frac{P_{Ik}\eta_k}{\lambda_1 \Lambda_k} \middle| \begin{matrix} 1 - m_{sk}, 1 - N_{Ik}, 1, 1 - \lambda_2 \\ m_{fk}, 0 \end{matrix} \right). \quad (16)$$

*Proof*: Please refer to appendix D. ■

##### B. Triplet Drop Probability

We consider that the BEP of  $k_{th}$  user is  $E_k$ . The TDP  $P_k$  can be expressed as

$$P_k = \sum_{j=D_E+1}^{D_T} E_k^j (1 - E_k)^{D_T-j}, \quad (17)$$

which can be calculated with the help of (15).

#### V. SEMANTIC LEVEL: PERSONALIZED SALIENCY FUSED SEMANTIC COMMUNICATION FRAMEWORK

In applications such as real-time sensing of UAV aerial photography, the autonomous patrol UAV cruises along the specified trajectory and transmits aerial pictures back to users (e.g., monitoring centers and photographers). In such a wireless communication environment with fading channels, transmitting all aerial images to all users is expensive, which is impractical in reality and limits the application of UAV aerial photography subscription. In this work, we propose that users have the right to specify which type of pictures they prefer, then the UAV only needs to transmit a small number of specific images to each user. For example, user A only needs to download pictures of “man wearing a jacket”, and user B prefers pictures of “bus on the street.” These texts are called personal queries, based on which the user

$$f_{\gamma_k}(\gamma) = \frac{\Lambda_k^{m_{sk}} (\sigma^2)^{N_{Ik} + m_{fk}} \gamma^{m_{fk} - 1}}{(P_{Ik}\eta_k)^{N_{Ik}} (\Lambda_k + \gamma\sigma^2)^{m_{fk} + m_{sk}} \Gamma(N_{Ik}) \Gamma(m_{sk}) \Gamma(m_{fk})} \times G_{0,0:0,2;1,1}^{1,0:2,0;1,1} \left( \begin{matrix} N_{Ik} \\ - \end{matrix} \middle| \begin{matrix} - \\ m_{sk} - N_{Ik}, 0 \end{matrix} \middle| \begin{matrix} 1 - m_{fk} - m_{sk} \\ 0, 1 - m_{sk} \end{matrix} \middle| \frac{\sigma^2}{P_{Ik}\eta_k}, \frac{-\Lambda_k}{\Lambda_k + \gamma\sigma^2} \right) \quad (9)$$

$$F_{\gamma_k}(\gamma) = \frac{\Lambda_k^{-m_{fk}} (\sigma^2)^{N_{Ik} + m_{fk}} \gamma^{m_{fk}}}{(P_{Ik}\eta_k)^{N_{Ik}} \Gamma(N_{Ik}) \Gamma(m_{sk}) \Gamma(m_{fk})} \times H_{2,0:2,0;2,0;2,1}^{0,2,0;2,0;1,1;1} \left( \begin{matrix} P_{Ik}\eta_k \sigma^{-2} \\ -1 \\ \sigma^{-2} \Lambda_k \gamma^{-1} \end{matrix} \middle| \begin{matrix} \Theta_1, \Theta_2 \\ - \end{matrix} \left[ \begin{matrix} (1 - m_{sk} + N_{Ik}, 1) \\ - \end{matrix} \right] \left[ \begin{matrix} (1, 1) \\ - \end{matrix} \right] \left[ \begin{matrix} (1, 1) \\ (m_{sk}, 1) \end{matrix} \right] \left[ \begin{matrix} (1, 1) \\ (1 + m_{fk}, 1) \\ (m_{sk}, 1) \end{matrix} \right] \right) \quad (10)$$

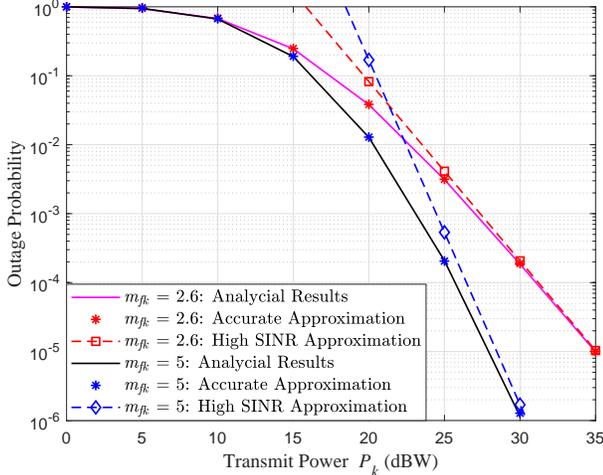


Fig. 2. The outage probability versus the transmit power, with different multipath fading parameter  $m_{fk}$ .

decides whether to download an image from the UAV. Does the problem come from *how do personal users know they need a given picture without downloading the original picture?* This paper addresses this problem by proposing a Personalized Saliency Fused Semantic Communication framework (PERSF-SEMCOM), which leverages a Triplet Detector to infer a fixed-size set of (subject-relation-object) triplets for each picture, and transmits these tiny-size triplets to users for triplet matching, in which only the matching images need to be downloaded. Besides, to prevent the triplets concerned by personalized users from being lost in the harsh wireless competitive environment, we quantify the priority of triplets and propose a transmission power allocation strategy, under which triplets with higher priority have more transmission power and therefore have a lower probability of being dropped.

#### A. OA-SemCom: A Fully Objective Approach

As shown in Figure 3, the architecture mainly contains Triplet Detection (TD), Personalized Saliency Prediction (PSP), Attention Fusion (AF), Triplet Priority Estimation (TPE), and Power Allocation (PA) on the UAV, and a Triplet Search (TS) module on the user terminal. Whenever the UAV captures a picture  $x$ , TD  $\mathcal{T}(x)$  infers a fixed-size set of (subject-relation-object) triplets, i.e., the semantic information,  $\tau$ , from  $x$ , and respectively the attention heatmaps  $H_{\text{sub}}, H_{\text{obj}}$  and box coordinates  $B_{\text{sub}}, B_{\text{obj}}$  of their subject and object entities. Specifically, the pretrained RelTR [38] model is used as the kernel of the TD module, which builds an encoder-decoder architecture like Transformer [45], where the encoder infers the visual feature context that is then used by the decoder for triplet inference. However, the TD module can also be implemented using any other pretrained scene graph generation technique. For example, two-stage approaches employ Fast/Faster R-CNN [46][47] to extract object features, and then apply scene graph generation [48][49] for graph inference. The alternative is one-stage approaches such as FCSGG [50] and the RelTR we use, which predict objects and their relations

concurrently, in an end-to-end fashion, and are thus more lightweight and faster. Regardless of which technique is used to implement TD, the general formula  $\mathcal{T}(x)$  for TD obeys

$$H_{\text{sub}}, H_{\text{obj}}, B_{\text{sub}}, B_{\text{obj}}, \tau \leftarrow \mathcal{T}(x), \quad (18)$$

where  $H_{\text{sub}}, H_{\text{obj}}$  are objective attention heatmaps of subjects and objects,  $B_{\text{sub}}, B_{\text{obj}}$  are bounding boxes of subjects and objects, respectively, and  $\tau$  is the prediction set of (subject-relation-object) triplets.

After this step, a naive idea would be to transmit the triplet set  $\tau$  to all users to match their personal queries  $q^k (\forall k \in [1, K])$ . This naive approach (named Naive-SemCom) faces the problem of key triplets being dropped due to intense competition among multiple users for scarce wireless channel resources. For example, the user  $k$  has the personal query  $q^k = \textit{man wearing jacket}$ , but unfortunately, the packet of the key triplet *man wearing jacket* is dropped in the wireless channel, then the current image will fail to match the personal query. The reason is that the wireless channel treats packets of all triplets as equally important, making the key triplets drop with equal probability as other triplets.

To this end, we propose prioritizing triplets for each user and allocating more transmission power to triplets with higher priority to ensure that key triplets are successfully delivered to the user terminal. To achieve this, one challenge should be addressed, that is, *how to prioritize triplets and identify the key triplets for personalized users?* As a benchmark, we can use the product of maximum values of the subject and object attention heatmaps  $H_{\text{sub}}, H_{\text{obj}}$  to obtain the triplet priority  $p$  in an objective view, namely,  $p = \max(\mathcal{C}(H_{\text{sub}}, B_{\text{sub}})) \otimes \max(\mathcal{C}(H_{\text{obj}}, B_{\text{obj}}))$ , where  $\mathcal{C}(H_{\text{sub}}, B_{\text{sub}})$  is a cropping function that crops the attention sub-image from  $H_{\text{sub}}$  according to the box coordinates  $B_{\text{sub}}$ . Note that the cropping function  $\mathcal{C}$  and  $\max$  in this section are channel-wise operations, and “ $\otimes$ ” is an element-wise multiplier. We refer to this benchmarking approach as Objective Attention-Based Semantic Communication (OA-SemCom), which, as the name suggests, only considers the objective global attention of the image itself, and its triplet priority is common to all users. However, users have personalized saliency, and their triplets should be prioritized differently than others, as are key triplets, but OA-SemCom fails to capture the personalization in subjective saliency among different users.

#### B. PERSF-SEMCOM: Achieve Personalization Through Fused Saliency

To address the above issues, we develop a PERSF-SEMCOM framework that integrates objective visual attention from RelTR and subjective visual attention from a personalized saliency prediction module, namely PSP, to assist the UAV in personalized priority estimation for each user. In our implementation, the PSP module  $\mathcal{P}(x, k)$  leverages a pretrained fully convolutional encoder-decoder network structure[51] to predict personal saliency heatmap  $S^k$  for user  $k$ ,

$$S^k \leftarrow \mathcal{P}(x, k). \quad (19)$$

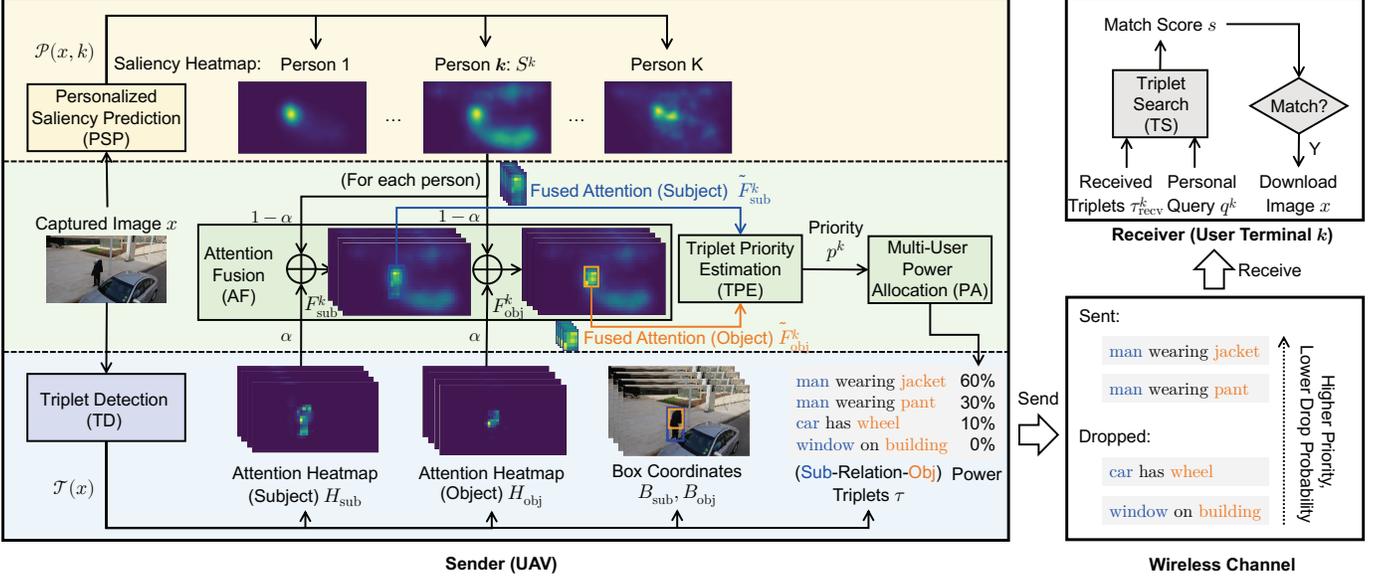


Fig. 3. An overview of the Personalized Saliency Fused Semantic Communication Framework (PERSF-SECOM).

These heatmaps show the things that different users are most interested in when looking at the same image: the users' subjective saliency distribution. For example, some users may only follow persons, while others also follow cars.

The encoder uses VGG16 [52] without pooling layers as the backbone, followed by an ASPP module to capture multi-scale visual information, and then the decoder restores the original image resolution by stacking convolution and up-sampling layers. Since how to achieve personalization in PSP is not the focus of this study, for simplicity, we use a plain but effective method to achieve personalization, which is to train user saliency models separately on different user datasets to obtain personalized behaviors. Please note that similar to the TD module, the PSP module is also a replaceable plugin that can be replaced with other pretrained personalized saliency models, depending on the researcher's preference. For instance, treat each user's prediction task separately and use a multi-task model to train a personalized saliency model for each user [33], or train a meta-learning model that can quickly adapt to new personalized tasks [53].

Then in Figure 3, we take a person  $k$  as an example to illustrate how to fuse objective attention and subjective saliency. In the AF module, the attention heatmaps  $H_{\text{sub}}, H_{\text{obj}}$  and the saliency heatmap  $S^k$  are first normalized, and then fused by weighted sum, where  $\alpha \in [0, 1]$  is the fusion coefficient, " $\oplus$ " is the broadcast add operator, and norm is a global normalizer,

$$F_{\text{sub}}^k = \alpha \cdot \text{norm}(S^k) \oplus (1 - \alpha) \cdot \text{norm}(H_{\text{sub}}), \quad (20)$$

$$F_{\text{obj}}^k = \alpha \cdot \text{norm}(S^k) \oplus (1 - \alpha) \cdot \text{norm}(H_{\text{obj}}). \quad (21)$$

In order to accurately locate the attention heatmaps of the subject and object entities, the fused attention heatmaps  $F_{\text{sub}}^k, F_{\text{obj}}^k$  should be cropped according to the box coordinates  $B_{\text{sub}}, B_{\text{obj}}$  to obtain the sub-heatmaps  $\tilde{F}_{\text{sub}}^k, \tilde{F}_{\text{obj}}^k$  of the subject and object

entities. These sub-heatmaps are then fed into the TPE module for priority estimation, where their maxima are multiplied and used as the triplet priorities  $p^k$ ,

$$\tilde{F}_{\text{sub}}^k = \mathcal{C}(F_{\text{sub}}^k, B_{\text{sub}}), \quad (22)$$

$$\tilde{F}_{\text{obj}}^k = \mathcal{C}(F_{\text{obj}}^k, B_{\text{obj}}), \quad (23)$$

$$p^k = \max(\tilde{F}_{\text{sub}}^k) \otimes \max(\tilde{F}_{\text{obj}}^k). \quad (24)$$

Finally, the PA module allocates transmission power for triplets according to personalized priorities  $p^k$ . The key triplets with higher priority will be allocated more transmission power, which makes them less likely to be dropped while traversing the wireless channel.

On the receiver side (i.e., the user terminal  $k$ ), the TS module uses the received triplets  $\tau_{\text{recv}}^k$  to match the personal query  $q^k$  and calculate a match score  $s$ . Here, we recommend two types of TS modules: Accurate Mode (TS-AM) and Fuzzy Mode (TS-FM). TS-AM aims to find the same received triplet as the personal query, and it returns a match score of 1 if found and 0 otherwise. TS-AM is preferred if the user's personal query is forced to meet the (subject-relation-object) format. However, if the personal query is free text, TS-FM could be better because it can return a match score at the semantic level (e.g., HEM[54]). Once a matching score  $s$  is obtained, the user can decide whether to download the current image  $x$ . In our implementation, we use TS-AM by default because semantic sentence matching is not the focus of this work, and only images with matching scores  $s = 1$  will be downloaded. We summarize the pseudo code of the proposed PERSF-SECOM in Algorithm 1.

**Algorithm 1** PERSF-SECOM (Main)**Input:** Captured image  $x$  on UAV, user identity  $k$ .**Output:** Match score  $s$ , downloaded image  $x$  on user  $k$ .

- 1: **procedure** UAV-SEND( $x, k$ )
- 2: TD detects triplets  $\tau$  and their attention heatmaps  $H_{\text{sub}}, H_{\text{obj}}$  and box coordinates  $B_{\text{sub}}, B_{\text{obj}}$ :  

$$H_{\text{sub}}, H_{\text{obj}}, B_{\text{sub}}, B_{\text{obj}}, \tau \leftarrow \mathcal{T}(x);$$
- 3: PSP predicts the personalized saliency heatmap  $S^k$  for user  $k$ :  $S^k \leftarrow \mathcal{P}(x, k)$ ;
- 4: AF fuses the objective attention  $H_{\text{sub}}, H_{\text{obj}}$  and the subjective saliency  $S^k$  by Eqs. (20)-(21);
- 5: TPE calculates the priority  $p^k$  of triplets  $\tau$  using the fused attention heatmaps  $F_{\text{sub}}^k, F_{\text{obj}}^k$  by Eqs. (22)-(24);
- 6: PA allocates transmission power to triplets  $\tau$  according to their priority  $p^k$  using RCGA[55];
- 7: UAV sends triplets  $\tau$  to user  $k$ 's terminal;
- 8: **procedure** USER-RECEIVE( $\tau_{\text{recv}}^k$ )
- 9: TS calculates the match score  $s$  between received triplets  $\tau_{\text{recv}}^k$  and personal query  $q^k$  via TS-AM/FM;
- 10: **if**  $s$  is greater than a user-specified threshold **then** User  $k$  downloads current image  $x$ ;  
**return** Match score  $s$  and downloaded image  $x$ ;

## VI. NUMERICAL RESULTS

## A. Environment Setup

The experimental platform is built on a generic Ubuntu 18.04 system with Intel(R) Xeon(R) E5-2678 CPU and 4 Geforce RTX 2080 TI GPUs. The RelTR model is adopted as the core of the TD module, which has been pretrained on the Visual Genome (VG) dataset [56], and has a top-50 recall rate of 25.2 on the scene graph detection metric, with the corresponding mean value being 8.5<sup>2</sup>. The VG dataset contains 108k images with 150 objects and 50 relationship categories. For the PSP module, we use the saliency prediction model in literature [51] pretrained on 3 visual attention datasets, respectively, including SALICON [57], MIT1003 [58] and DUT-OMRON [59], to simulate the visual saliency discrepancy of 3 users<sup>3</sup>. The reason this works is that researchers collect these datasets from different perspectives, which can be regarded as real-world experiences and environments of different persons, thus presenting independent personalities. The validation dataset used in our experiments is real-world video frames downloaded from YouTube<sup>4</sup>. We framed this video at an interval of 5 frames and obtained a dataset of 59 images, which we named ‘‘STREET’’. In Figure 4, we visualize a portion of the STREET dataset. Each user’s personal queries and their experience datasets are summarized in Table III.

If not otherwise specified, the parameters of the small-scale fading model and the system parameters such as power and transmission distance are shown in Table II following common parameter settings in the literature [60], [61].

<sup>2</sup>The pretrained model is available at: <https://github.com/yrcong/RelTR>

<sup>3</sup>The pretrained model is available at: <https://github.com/alexanderkroner/saliency>

<sup>4</sup>The video is available at: <https://www.youtube.com/watch?v=RPZ3xWy70IE>



Fig. 4. A partial overview of our STREET dataset.

TABLE II  
CHANNEL CONFIGURATION PARAMETERS FOR 3 USERS.

User ID	1	2	3
Fading parameter $m_f$	2	2	5
Shadowing parameter $m_s$	2	4	2
Signal amplitude decrease $\bar{z}$ (dB)	-3		
Distance $D_{uk}$ (m)	10		
Number of antennas $N_T$	3		
Paths of interferes $N_{I,k}$	2		
Interference power $P_{I,k,j}$ (W)	2		
$\eta_k$ (dB)	-3		
Noise $n_k$ (W)	1		
$\tau_1$	1		
$\tau_2$	0.5		
Path loss exponent $\alpha_k$	2		

Two benchmark approaches are used for performance comparison, that is, Naive-SemCom and OA-SemCom. Naive-SemCom aims to allocate transmit power when transmitting to each user equally, and the triplets of each user are also allocated with equal power. In other words, Naive-SemCom has no awareness of the importance of semantic triples. Instead, OA-SemCom quantifies the priority of triplets and uses these priorities to schedule transmit power for each user and triplet. However, OA-SemCom only considers the objective attention obtained by TD but ignores the user’s subjective attention and personality. To address this issue, our PERSF-SECOM is proposed, and PSP, AF modules are introduced to fuse the objective and subjective attention. We set the fusion coefficient  $\alpha = 0.2$  and the transmit power  $P = 3000$  by default and calculated the average of each user’s match score on all images as the user’s score. Each experiment was repeated 5 times, and the average results are shown.

## B. Results and Analysis

**Effectiveness of PERSF-SECOM over  $\alpha$  and  $P$ .** We first show the effectiveness of the proposed PERSF-SECOM. To find the optimal hyperparameters, we increase the fusion coefficient  $\alpha$  from 0 to 1 and the total transmits power  $P$  from 1kW to 3kW. The utility value (i.e., the product of

TABLE III  
USERS' IDENTITIES, EXPERIENCE DATASETS AND PERSONAL QUERIES.

Identity	Dataset	Personal Queries
User 1	SALICON	woman has hair
User 2	MIT1003	sign on building
User 3	DUT-OMRON	woman wearing shirt

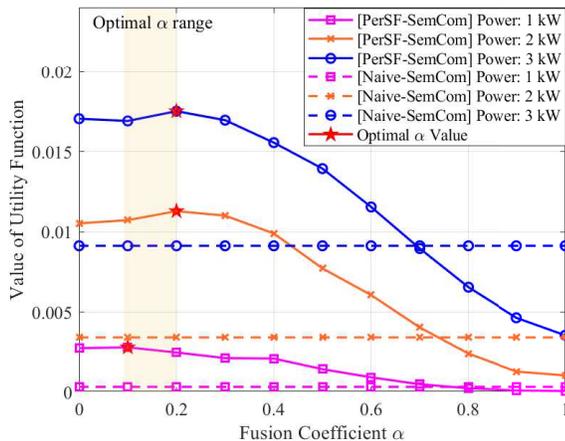


Fig. 5. The value curves of utility function over different fusion coefficient  $\alpha$  and transmit power  $P$ .

all user scores) curves are shown in Figure 5. Please note that OA-SemCom corresponds to our PERSF-SEMCOM with  $\alpha = 1.0$ . Benefiting from the introduction of subjective attention, PERSF-SEMCOM outperforms Naive-SemCom in most cases, for example, when  $\alpha \in [0, 0.6]$  and  $P \in [1\text{kW}, 3\text{kW}]$ . The curves of PERSF-SEMCOM first increase to the optimal and then rapidly degrade. The optimal utility value usually occurs when  $\alpha$  is between 0.1 and 0.2; that is, the objective attention should contribute 10%-20% to the fused attention, while the subjective attention contributes 80%-90%. After this, the PERSF-SEMCOM performance gradually deteriorated with the withdrawal of subjective attention, especially when  $\alpha = 1.0$ , PERSF-SEMCOM degenerates to OA-SemCom, which performs even worse than Naive-SemCom. The above results show that the introduction of appropriate subjective attention can significantly enhance semantic communication's personalization and anti-interference ability and demonstrates the effectiveness of the proposed PERSF-SEMCOM. In subsequent experiments, we use the recommended  $\alpha = 0.2$  as the default setting.

**The effects of transmit power  $P$ .** To explore the optimality gap over different transmit powers  $P$ , we increase  $P$  from 1kW to 3kW and illustrate the gap between PERSF-SEMCOM and the theoretical optimal in Figure 6. The top dashed lines represent the theoretical upper bounds of each user's score. They use infinite transmit power, so no triplet packets get dropped. Users have different upper bounds because of their personality divergence. As expected, the optimality gap of PERSF-SEMCOM becomes smaller as the transmit power  $P$  increases because higher transmit power effectively reduces the probability of packet loss. Compared with Naive-SemCom,

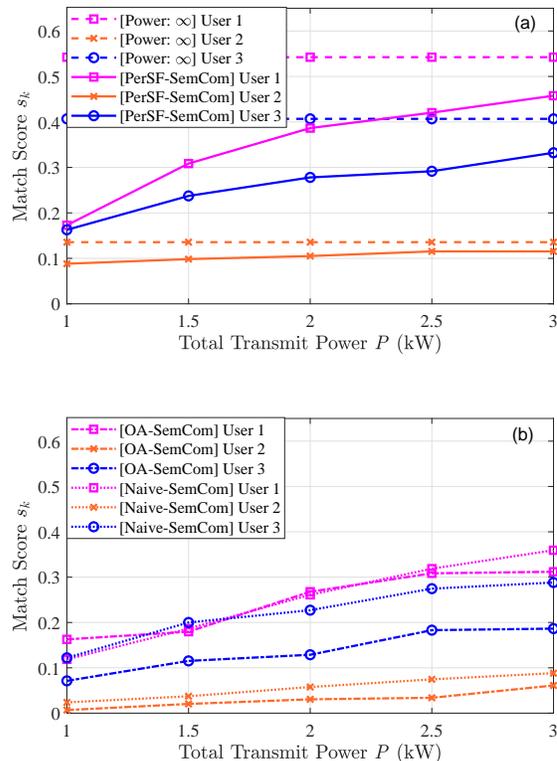


Fig. 6. The score curves for each user over different total transmit power  $P$ .

when  $P = 3\text{kW}$ , PERSF-SEMCOM reduces the optimality gap by 54%, 57%, 37% on the three users, respectively, which shows a significant improvement in the accuracy of the UAV aerial imagery subscription service. However, OA-SemCom, which only considers objective attention, widens the optimality gap by 26%, 57%, and 86%. These results prove the necessity to incorporate personal attention.

**The effects of power allocation.** Intuitively, users with poor channel conditions should be allocated more transmission power. In this experiment, the channel condition ranking of users 1-3 is user 1 < user 3 < user 2. We use the Real-Coded Genetic Algorithm (RCGA)[55] to solve the NBS problem and allocate power among users, and then proportionally distribute power among triplets according to priority. RCGA uses the population size of 50, the mutation probability of 0.001, and the maximum iteration of 20. Given the transmit power  $P = 3\text{kW}$  and 3 users, we traverse all possible settings of full power allocation (i.e., the sum of the proportions of the power allocated to each user is 1) and visualize their utility function surface in Figure 7. RCGA found the best power allocation per user to be (40.6%, 18.4%, 40.0%), with the score (0.46, 0.12, 0.33) per user and the utility value 0.0175, which outperforms Naive-SemCom with the power allocation (33.3%, 33.3%, 33.3%), the scores (0.36, 0.09, 0.29) and the utility value 0.0091. The resulting power allocation strategy also supports our intuitive idea that more power should be allocated to users with poor channel conditions rather than

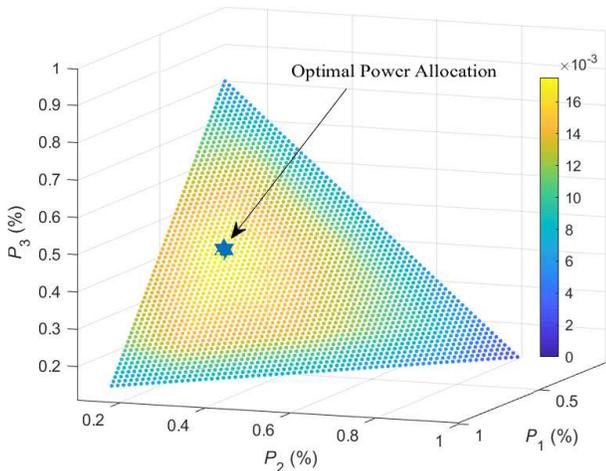


Fig. 7. The utility function surface under different power allocations.

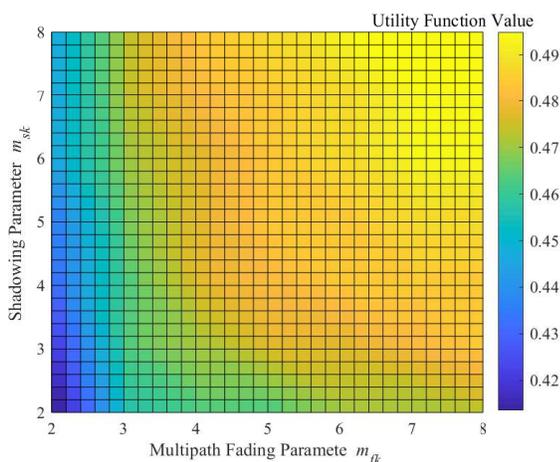


Fig. 8. Utility function values under different multi-path fading parameter  $m_{fk}$  and shadow fading parameter  $m_{sk}$ .

simply evenly allocated.

**The effects of small-scale channel conditions  $m_{sk}, m_{fk}$ .** Given the transmission power of user  $k$  is 1kW, we investigate the effect of different values of multi-path fading parameter  $m_{fk}$  and shading parameter  $m_{sk}$  on the utility function values. As shown in Fig. 8, when the values of both  $m_{fk}$  and  $m_{sk}$  are large, which means the multipath effect is weak and there is less shading, the value of the utility function is large. A more interesting insight is that an increase in  $m_{fk}$  leads to a faster increase in the utility function value compared to the same increase in  $m_{sk}$ . This suggests that, at the wireless channel level, the goal-oriented semantic communication system that we studied in this paper is more affected by the BER increase due to the multi-path effect, instead of shadowing.

**The effects of large-scale channel conditions  $D_{u,k}, P_{Ik}$ .** Given the transmission power of user  $k$  is 1kW, we investigate the effect of transmission distance and interference power on the value of the utility function. In Fig. 9, we can observe that, when the interference power is small, i.e.,  $P_{Ik} < 1$  W, the increase in the transmission distance does not result in a sig-

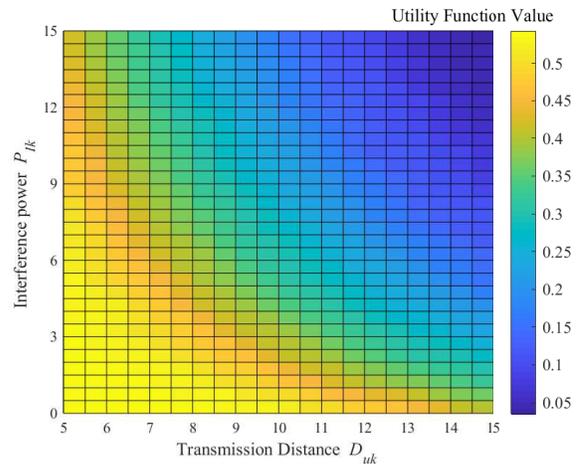


Fig. 9. Utility function values under different transmission distance  $D_{uk}$  and interference power  $P_{Ik}$ .

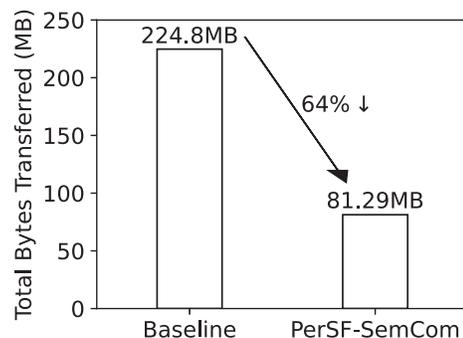


Fig. 10. Comparison of total transferred bytes of images.

nificant decrease in the value of the utility function. However, in the high interference regime, e.g., when  $P_{Ik} > 10$  W, every 5 m increase in transmission distance results in an about 58% decrease in the value of the utility function. Therefore, when there is substantial interference in the environment, we need to reduce the transmission distance by adjusting the trajectory of the UAV to ensure the quality of the semantic communication services.

**The effects of reducing communication overhead.** Considering a vanilla approach that the UAV sends a total of 59 images to 3 subscribers, each image is 1.27MB in size and the total bytes transferred is 224.8MB. As an improvement, PERSF-SEMCOM selectively sends fewer images (64 in total) to specified subscribers, before that 873 triplets were sent with a negligible additional communication cost of 10.24KB. Then, the UAV only needs to send 81.29MB in size, which reduces the communication cost by 64% and thus saves the power consumption during transmission.

## VII. CONCLUSIONS

In this paper, we have focused on the semantic communication personalization and resource allocation optimization issues for personalized saliency-based task-oriented semantic communication in UAV image sensing scenarios. We have first

presented an energy-efficient task-oriented semantic communication framework with an efficient image retrieval manner based on a triple-based scene graph. To ensure personalized semantic communication, we have designed a personalized attention-based mechanism to realize differential weight encoding of triplets for important information according to user preferences. Furthermore, we have analyzed mathematically the effects of wireless fading channels on semantic communication and proposed a game-based model for a multi-user resource allocation scheme to achieve efficient utilization of UAV resources. We evaluate performance of the proposed framework and schemes on real-world datasets. The numerical results have confirmed that the proposed framework and schemes can realize personalized semantic communication and significantly enhance the UAV resource utilization.

## APPENDIX A PROOF OF LEMMA 1

### A. Proof of PDF

Let  $X \triangleq \sigma^2 + P_{Ik}Y$  and  $U \triangleq P_k D_{uk}^{-\alpha_k} Z$ . Thus, we have  $\gamma_k = \frac{U}{X}$ . The PDF of  $\gamma_k$  can be expressed as

$$f_{\gamma_k}(\gamma) = \int_0^\infty x f_U(\gamma x) f_X(x) dx. \quad (\text{A-1})$$

Substituting (7) and [41, eq. (6)] into (A-1), we have

$$f_{\gamma_k}(\gamma) = \frac{m_{fk} m_{fk} (m_{sk} 1)^{m_{sk}} \bar{z}_k^{m_{sk}} \gamma^{m_{fk} 1}}{(P_{Ik} \eta_k)^{N_{Ik}} \Gamma(K) (P_k D_{uk}^{\alpha_k})^{m_{fk}} B(m_{fk}, m_{sk})} I_1, \quad (\text{A-2})$$

where

$$I_1 = \int_{\sigma^2}^\infty \frac{x^{m_{fk}} (x\sigma^2)^{N_{Ik} 1} \exp\left(\frac{x\sigma^2}{P_{Ik} \eta_k}\right)}{\left(\frac{m_{fk} \gamma x}{P_k D_{uk}^{\alpha_k}} + (m_{sk} 1) \bar{z}_k\right)^{m_{fk} + m_{sk}}} dx. \quad (\text{A-3})$$

With the help of [20, eq. (9.301)] and [43, eq. (01.03.26.0004.01)], we can express the  $\exp(\cdot)$  function in terms of the Mellin-Barnes integral form as

$$\begin{aligned} \exp\left(-\frac{x - \sigma^2}{P_{Ik} \eta_k}\right) &= G_{0,1}^{1,0}\left(\frac{x - \sigma^2}{P_{Ik} \eta_k} \middle| 0\right) \\ &= \frac{1}{2\pi i} \int_{\mathcal{L}_1} \Gamma(-s_1) \left(\frac{x - \sigma^2}{P_{Ik} \eta_k}\right)^{s_1} ds_1. \end{aligned} \quad (\text{A-4})$$

Substituting (A-4) into  $I_1$ , we obtain that

$$\begin{aligned} I_1 &= \frac{1}{2\pi i} \int_{\mathcal{L}_1} \Gamma(-s_1) \left(\frac{1}{P_{Ik} \eta_k}\right)^{s_1} \\ &\times \int_0^\infty \frac{x^{s_1 + N_{Ik} - 1} (x + \sigma^2)^{m_{fk}} dx ds_1}{\left(\frac{m_{fk} \gamma x}{P_k D_{uk}^{\alpha_k}} + \frac{m_{fk} \gamma \sigma^2}{P_k D_{uk}^{\alpha_k}} + (m_{sk} - 1) \bar{z}_k\right)^{m_{fk} + m_{sk}}}. \end{aligned} \quad (\text{A-5})$$

Let  $\Lambda_k \triangleq \frac{P_k D_{uk}^{-\alpha_k} (m_{sk} - 1) \bar{z}_k}{m_{fk}}$ . With the help of [20, eq. (3.197.1)], the integration part in  $I_1$  can be solved. Thus, we

can re-write  $I_1$  as

$$\begin{aligned} I_1 &= \frac{1}{2\pi i} \left(\frac{P_k D_{uk}^{-\alpha_k}}{m_{fk} \gamma}\right)^{m_{fk} + m_{sk}} \int_{\mathcal{L}} \Gamma(-s_1) \\ &\times \left(\frac{1}{P_{Ik} \eta_k}\right)^{s_1} \left(\frac{P_k D_{uk}^{-\alpha_k} (m_{sk} - 1) \bar{z}_k}{m_{fk} \gamma} + \sigma^2\right)^{-m_{fk} - m_{sk}} \\ &\times \sigma^{2s_1 + 2N_{Ik} + 2m_{fk}} B(s_1 + N_{Ik}, m_{sk} - s_1 - N_{Ik}) \\ &\times F\left(m_{fk} + m_{sk}, s_1 + N_{Ik}; m_{sk}; 1 - \frac{\sigma^2}{\frac{\Lambda_k}{\gamma} + \sigma^2}\right) ds_1. \end{aligned} \quad (\text{A-6})$$

Using [20, eq. (9.113)] and [20, eq. (8.384.1)], we can further express  $I_1$  as

$$\begin{aligned} I_1 &= \left(\frac{1}{2\pi i}\right)^2 \left(\frac{P_k D_{uk}^{-\alpha_k}}{m_{fk} (\Lambda_k + \gamma \sigma^2)}\right)^{m_{fk} + m_{sk}} \frac{(\sigma^2)^{N_{Ik} + m_{fk}}}{\Gamma(m_{fk} + m_{sk})} \\ &\times \int_{\mathcal{L}_1} \int_{\mathcal{L}_2} \frac{\Gamma(m_{sk} - s_1 - N_{Ik}) \Gamma(-s_1) \Gamma(m_{fk} + m_{sk} + s_2)}{\Gamma(m_{sk} + s_2) \Gamma^{-1}(-s_2) \Gamma^{-1}(s_1 + N_{Ik} + s_2)} \\ &\times \left(\frac{\sigma^2}{P_{Ik} \eta_k}\right)^{s_1} \left(\frac{-\Lambda_k}{\Lambda_k + \gamma \sigma^2}\right)^{s_2} ds_2 ds_1. \end{aligned} \quad (\text{A-7})$$

Therefore, substituting  $I_1$  into (A-2), using the definition of Bivariate Meijer's  $G$ -function [22, eq. (1)], we can derive (9) to complete the proof.

### B. Proof of CDF

According to the definition of CDF, we have

$$F_{\gamma_k}(\gamma) = \int_0^\gamma f_{\gamma_k}(x) dx. \quad (\text{A-8})$$

Combining (9) with (A-8), we have

$$\begin{aligned} F_{\gamma_k}(\gamma) &= \frac{m_{fk} m_{fk} (m_{sk} - 1)^{m_{sk}} \bar{z}_k^{m_{sk}}}{(P_{Ik} \eta_k)^{N_{Ik}} \Gamma(N_{Ik}) (P_k D_{uk}^{-\alpha_k})^{m_{fk}} B(m_{fk}, m_{sk})} \\ &\times \left(\frac{1}{2\pi i}\right)^2 \left(\frac{P_k D_{uk}^{-\alpha_k}}{m_{fk}}\right)^{m_{fk} + m_{sk}} \frac{(\sigma^2)^{N_{Ik} + m_{fk}}}{\Gamma(m_{fk} + m_{sk})} \\ &\times \int_{\mathcal{L}_1} \int_{\mathcal{L}_2} \frac{\Gamma(m_{sk} - s_1 - N_{Ik}) \Gamma(-s_1) \Gamma(m_{fk} + m_{sk} + s_2)}{\Gamma(m_{sk} + s_2) \Gamma^{-1}(-s_2) \Gamma^{-1}(s_1 + N_{Ik} + s_2)} \\ &\times I_2 \left(\frac{\sigma^2}{P_{Ik} \eta_k}\right)^{s_1} (-\Lambda_k)^{s_2} ds_2 ds_1, \end{aligned} \quad (\text{A-9})$$

where the  $I_2$  can be expressed as

$$I_2 = \int_0^\gamma \frac{x^{m_{fk} - 1}}{(\Lambda_k + x\sigma^2)^{m_{fk} + m_{sk} + s_2}} dx. \quad (\text{A-10})$$

With the help of, we can solve  $I_2$  as

$$\begin{aligned} I_2 &= \frac{\gamma^{m_{fk}}}{\Lambda_k^{m_{fk} + m_{sk} + s_2} m_{fk}} \\ &\times F\left(m_{fk} + m_{sk} + s_2, m_{fk}; 1 + m_{fk}; -\frac{\sigma^2}{\Lambda_k} \gamma\right). \end{aligned} \quad (\text{A-11})$$

According to the integral expression of the hyper-geometric function [20, eq. (9.113)], we can substitute  $I_2$  into (A-9) and obtain (10), which completes the proof.

APPENDIX B  
PROOF OF LEMMA 1

Let  $t \triangleq x\rho$ . We can re-write  $I_A$  as

$$I_A = \rho^{-b-c-1} \exp\left(\frac{a}{d}\right) \times \int_{\rho a}^{\infty} t^b (t - \rho a)^c \exp\left(-\frac{t}{\rho d}\right) F(\alpha, \beta, \varepsilon; -t) dt. \quad (\text{B-1})$$

Because  $\rho$  is small, we can further express  $I_A$  as

$$I_A \approx \rho^{-b-c-1} \exp\left(\frac{a}{d}\right) \times \int_0^{\infty} t^{b+c} \exp\left(-\frac{t}{\rho d}\right) F(\alpha, \beta, \varepsilon; -t) dt. \quad (\text{B-2})$$

With the help of [20, eq. (9.113)], we obtain

$$I_A = \rho^{-b-c-1} \exp\left(\frac{a}{d}\right) \frac{\Gamma(\varepsilon)}{\Gamma(\alpha)\Gamma(\beta)} \frac{1}{2\pi i} \times \int_{\mathcal{L}_1} \frac{\Gamma(s_1 + \alpha)\Gamma(s_1 + \beta)\Gamma(-s_1)}{\Gamma(s_1 + \varepsilon)} I_3 ds_1, \quad (\text{B-3})$$

where

$$I_3 = \int_0^{\infty} t^{s_1+b+c} \exp\left(-\frac{t}{\rho d}\right) dt. \quad (\text{B-4})$$

By using [20, eq. (3.351.3)] and [20, eq. (8.339.1)],  $I_3$  can be solved as

$$I_3 = (d\rho)^{b+c+s_1+1} \Gamma(1 + b + c + s_1). \quad (\text{B-5})$$

Substituting  $I_3$  into (B-3), we have

$$I_A = \exp\left(\frac{a}{d}\right) \frac{\Gamma(\varepsilon) d^{b+c+1}}{\Gamma(\alpha)\Gamma(\beta)} \frac{1}{2\pi i} \times \int_{\mathcal{L}_1} \frac{\Gamma(s_1 + \alpha)\Gamma(s_1 + \beta)\Gamma(-s_1)}{\Gamma^{-1}(1 + b + c + s_1)\Gamma(s_1 + \varepsilon)} (d\rho)^{s_1} ds_1. \quad (\text{B-6})$$

According to the definition of Meijer's  $G$ -function [20, eq. (9.301)], we can re-write  $I_A$  as (12) to complete the proof.

APPENDIX C  
PROOF OF THEOREM 2

Using the definition of CDF, we have

$$F_{\gamma_k}(\gamma) = \int_0^{\infty} F_U(\gamma x) f_X(x) dx, \quad (\text{C-1})$$

where

$$F_U(u) = \Pr(U < u) = \Pr\left(Z < \frac{u}{P_k D_{uk}^{-\alpha_k}}\right) = F_Z\left(\frac{u}{P_k D_{uk}^{-\alpha_k}}\right). \quad (\text{C-2})$$

Thus, the CDF of  $\gamma_k$  can be expressed as

$$F_{\gamma_k}(\gamma) = \frac{1}{m_{fk} B(m_{fk}, m_{sk})} \left(\frac{\gamma}{\Lambda_k}\right)^{m_{fk}} \frac{1}{(P_{Ik}\eta_k)^{N_{Ik}} \Gamma(N_{Ik})} \times \int_{\sigma^2}^{\infty} x^{m_{fk}} (x - \sigma^2)^{N_{Ik}-1} \exp\left(-\frac{x - \sigma^2}{P_{Ik}\eta_k}\right) \times F\left(m_{fk}, m_{fk} + m_{sk}, m_{fk} + 1; -\frac{\gamma x}{\Lambda_k}\right) dx. \quad (\text{C-3})$$

The integration part in (C-3) can be solved with the help of Lemma 1. Then, after some algebraic manipulations, we have

$$F_{\gamma_k}(\gamma) = \frac{1}{\Gamma(N_{Ik}) \Gamma(m_{sk}) \Gamma(m_{fk})} \exp\left(\frac{\sigma^2}{P_{Ik}\eta_k}\right) \times \frac{1}{2\pi i} \int_{\mathcal{L}} \frac{\Gamma(s_2) \Gamma(s_2 + m_{sk}) \Gamma(N_{Ik} + s_2)}{\Gamma(s_2 + 1) \Gamma^{-1}(-s_2 + m_{fk})} \left(\frac{\gamma P_{Ik}\eta_k}{\Lambda_k}\right)^{s_2} ds_2. \quad (\text{C-4})$$

Using [20, eq. (9.301)], we can derive the CDF of  $\gamma_k$  as (10), which completes the proof.

APPENDIX D  
PROOF OF THEOREM 4

Using the definition of Gamma function [20, eq. (8.350)], we can express  $E_k$  as

$$E_k = \frac{\lambda_1 \lambda_2}{2\Gamma(\lambda_2)} \int_0^{\infty} x^{\lambda_2-1} e^{-\lambda_1 x} F_{\gamma_k}(x) dx. \quad (\text{D-1})$$

By substituting (C-4) into (D-1), we obtain

$$E_k = \frac{\lambda_1 \lambda_2 \Gamma^{-1}(m_{sk}) \Gamma^{-1}(m_{fk})}{2\Gamma(\lambda_2) \Gamma(N_{Ik})} \exp\left(\frac{\sigma^2}{P_{Ik}\eta_k}\right) \frac{1}{2\pi i} \times \int_{\mathcal{L}} \frac{\Gamma(s_2) \Gamma(s_2 + m_{sk}) \Gamma(N_{Ik} + s_2)}{\Gamma(s_2 + 1) \Gamma^{-1}(-s_2 + m_{fk})} \left(\frac{P_{Ik}\eta_k}{\Lambda_k}\right)^{s_2} I_4 ds_2, \quad (\text{D-2})$$

where

$$I_4 = \int_0^{\infty} x^{s_2+\lambda_2-1} e^{-\lambda_1 x} dx. \quad (\text{D-3})$$

Using [20, eq. (3.351.3)], we can solve  $I_4$  as

$$I_4 = \lambda_1^{-s_2-\lambda_2} \Gamma(s_2 + \lambda_2). \quad (\text{D-4})$$

Combining  $I_4$  and (D-2), we derive  $E_k$  as (16) to complete the proof.

REFERENCES

- [1] S. Dang, O. Amin, B. Shihada, and M.-S. Alouini, "What should 6g be?" *Nat. Electron.*, vol. 3, no. 1, pp. 20–29, 2020.
- [2] Y. Liu, X. Yuan, Z. Xiong, J. Kang, X. Wang, and D. Niyato, "Federated learning for 6g communications: Challenges, methods, and future directions," *China Commun.*, vol. 17, no. 9, pp. 105–118, Sept. 2020.
- [3] M. Giordani, M. Polese, M. Mezzavilla, S. Rangan, and M. Zorzi, "Toward 6g networks: Use cases and technologies," *IEEE Commun. Mag.*, vol. 58, no. 3, pp. 55–61, 2020.
- [4] X. You, C.-X. Wang, J. Huang, X. Gao, Z. Zhang, M. Wang, Y. Huang, C. Zhang, Y. Jiang, J. Wang *et al.*, "Towards 6g wireless communication networks: Vision, enabling technologies, and new paradigm shifts," *Science China Information Sciences*, vol. 64, no. 1, pp. 1–74, 2021.
- [5] L. Zhang, Y.-C. Liang, and D. Niyato, "6g visions: Mobile ultra-broadband, super internet-of-things, and artificial intelligence," *China Commun.*, vol. 16, no. 8, pp. 1–14, Aug 2019.
- [6] W. Yang, H. Du, Z. Liew, W. Y. B. Lim, Z. Xiong, D. Niyato, X. Chi, X. S. Shen, and C. Miao, "Semantic communications for 6G future internet: Fundamentals, applications, and challenges," *arXiv preprint arXiv:2207.00427*, 2022.
- [7] W. Saad, M. Bennis, and M. Chen, "A vision of 6G wireless systems: Applications, trends, technologies, and open research problems," *IEEE Netw.*, vol. 34, no. 3, pp. 134–142, Mar. 2019.
- [8] M. Xu, W. C. Ng, W. Y. B. Lim, J. Kang, Z. Xiong, D. Niyato, Q. Yang, X. S. Shen, and C. Miao, "A full dive into realizing the edge-enabled metaverse: Visions, enabling technologies, and challenges," *arXiv preprint arXiv:2203.05471*, 2022.

- [9] P. Zhang, W. Xu, H. Gao, K. Niu, X. Xu, X. Qin, C. Yuan, Z. Qin, H. Zhao, J. Wei *et al.*, "Toward wisdom-evolutionary and primitive-concise 6g: A new paradigm of semantic communication networks," *Engineering*, vol. 8, pp. 60–73, 2022.
- [10] E. C. Strinati and S. Barbarossa, "6g networks: Beyond shannon towards semantic and goal-oriented communications," *Computer Netw.*, vol. 190, p. 107930, 2021.
- [11] Q. Lan, D. Wen, Z. Zhang, Q. Zeng, X. Chen, P. Popovski, and K. Huang, "What is semantic communication? a view on conveying meaning in the era of machine intelligence," *J. Commun. Netw.*, vol. 6, no. 4, pp. 336–371, 2021.
- [12] G. Shi, Y. Xiao, Y. Li, and X. Xie, "From semantic communication to semantic-aware networking: Model, architecture, and open problems," *IEEE Commun. Mag.*, vol. 59, no. 8, pp. 44–50, Aug. 2021.
- [13] H. Xie, Z. Qin, and G. Y. Li, "Task-oriented multi-user semantic communications for vqa task," *IEEE Wireless Commun. Lett.*, 2021.
- [14] Q. Zhou, R. Li, Z. Zhao, C. Peng, and H. Zhang, "Semantic communication with adaptive universal transformer," *IEEE Wireless Commun. Lett.*, 2021.
- [15] X. Kang, B. Song, J. Guo, Z. Qin, and F. R. Yu, "Task-oriented image transmission for scene classification in unmanned aerial systems," *arXiv preprint arXiv:2112.10948*, 2021.
- [16] O. Ayan, P. Kutsevol, H. Y. Özkan, and W. Kellerer, "Task-oriented scheduling for networked control systems: An age of information-aware implementation on software-defined radios," *arXiv preprint arXiv:2202.09189*, 2022.
- [17] H. Xie and Z. Qin, "A lite distributed semantic communication system for internet of things," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 1, pp. 142–153, Jan. 2020.
- [18] O. Friha, M. A. Ferrag, L. Shu, L. A. Maglaras, and X. Wang, "Internet of things for the future of smart agriculture: A comprehensive survey of emerging technologies," *IEEE CAA J. Autom. Sinica*, vol. 8, no. 4, pp. 718–752, Apr. 2021.
- [19] H. Xie, Z. Qin, G. Y. Li, and B.-H. Juang, "Deep learning enabled semantic communication systems," *IEEE Trans. Signal Process.*, vol. 69, pp. 2663–2675, 2021.
- [20] I. S. Gradshteyn and I. M. Ryzhik, *Table of Integrals, Series, and Products*, 7th ed. Academic Press, 2007.
- [21] A. M. Mathai, R. K. Saxena, and H. J. Haubold, *The H-function: Theory and Applications*. Springer Science & Business Media, 2009.
- [22] B. L. Sharma and R. F. A. Abiodun, "Generating function for generalized function of two variables," *Proc. American Mathematical Society*, vol. 46, no. 1, pp. 69–72, Oct. 1974.
- [23] Z. Qin, X. Tao, J. Lu, and G. Y. Li, "Semantic communications: Principles and challenges," *arXiv preprint arXiv:2201.01389*, 2021.
- [24] E. Boursoulatzé, D. B. Kurka, and D. Gündüz, "Deep joint source-channel coding for wireless image transmission," *IEEE Trans. Cogn. Develop. Syst.*, vol. 5, no. 3, pp. 567–579, 2019.
- [25] M. K. Farshbafan, W. Saad, and M. Debbah, "Common language for goal-oriented semantic communications: A curriculum learning framework," *arXiv preprint arXiv:2111.08051*, 2021.
- [26] M. Chen, Z. Yang, W. Saad, C. Yin, H. V. Poor, and S. Cui, "A joint learning and communications framework for federated learning over wireless networks," *IEEE Trans. Wireless Commun.*, vol. 20, no. 1, pp. 269–283, Jan. 2020.
- [27] Z. Yang, M. Chen, W. Saad, C. S. Hong, and M. Shikh-Bahaee, "Energy efficient federated learning over wireless communication networks," *IEEE Trans. Wireless Commun.*, vol. 20, no. 3, pp. 1935–1949, Mar. 2020.
- [28] H. Tong, Z. Yang, S. Wang, Y. Hu, W. Saad, and C. Yin, "Federated learning based audio semantic communication over wireless networks," in *Proc. of GLOBECOM*. IEEE, 2021.
- [29] S. K. Yoo, S. L. Cotton, P. C. Sofotasios, M. Matthaiou, M. Valkama, and G. K. Karagiannidis, "The Fisher–snedecor  $\mathcal{F}$  distribution: A simple and accurate composite fading model," *IEEE Commun. Lett.*, vol. 21, no. 7, pp. 1661–1664, Jul. 2017.
- [30] S. K. Yoo, P. C. Sofotasios, S. L. Cotton, S. Muhaidat, F. J. Lopez-Martinez, J. M. Romero-Jerez, and G. K. Karagiannidis, "A comprehensive analysis of the achievable channel capacity in  $\mathcal{F}$  composite fading channels," *IEEE Access*, vol. 7, pp. 34 078–34 094, 2019.
- [31] H. Xie and Z. Qin, "A lite distributed semantic communication system for internet of things," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 1, pp. 142–153, 2021.
- [32] Z. Weng and Z. Qin, "Semantic communication systems for speech transmission," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 8, pp. 2434–2444, 2021.
- [33] Y. Xu, S. Gao, J. Wu, N. Li, and J. Yu, "Personalized saliency and its prediction," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 12, pp. 2975–2989, Dec. 2018.
- [34] S. F. Dodge and L. J. Karam, "Visual saliency prediction using a mixture of deep neural networks," *IEEE Trans. Image Process.*, vol. 27, no. 8, pp. 4080–4090, Aug. 2018.
- [35] A. Mahdi, J. Qin, and G. Crosby, "Deepfeat: A bottom-up and top-down saliency model based on deep features of convolutional neural networks," *IEEE Trans. Cogn. Develop. Syst.*, vol. 12, no. 1, pp. 54–63, 2019.
- [36] S. Berkovsky, R. Taib, I. Koprinska, E. Wang, Y. Zeng, J. Li, and S. Kleitman, "Detecting personality traits using eye-tracking data," in *Proc. CHI Conf. Hum. Factors Comput. Syst.*, 2019, pp. 1–12.
- [37] Y. Moroto, K. Maeda, T. Ogawa, and M. Haseyama, "Few-shot personalized saliency prediction using person similarity based on collaborative multi-output gaussian process regression," in *Proc. of ICIP*. IEEE, 2021, pp. 1469–1473.
- [38] Y. Cong, M. Y. Yang, and B. Rosenhahn, "Reltr: Relation transformer for scene graph generation," *arXiv preprint arXiv:2201.11460*, 2022.
- [39] T. K. Lo, "Maximum ratio transmission," in *IEEE Int. Conf. Commun. (Cat. No. 99CH36311)*, vol. 2, 1999, pp. 1310–1314.
- [40] H. Du, J. Zhang, J. Cheng, and B. Ai, "Sum of Fisher–snedecor  $\mathcal{F}$  random variables and its applications," *IEEE Open J. Commun. Soc.*, vol. 1, pp. 342–356, Mar. 2020.
- [41] S. K. Yoo, P. C. Sofotasios, S. L. Cotton, S. Muhaidat, F. J. Lopez-Martinez, J. M. Romero-Jerez, and G. K. Karagiannidis, "A comprehensive analysis of the achievable channel capacity in  $\mathcal{F}$  composite fading channels," *IEEE Access*, vol. 7, pp. 34 078–34 094, 2019.
- [42] M. Nakagami, "The m-distribution—a general formula of intensity distribution of rapid fading," in *Stat. methods Radio Wave Propag.* Elsevier, 1960, pp. 3–36.
- [43] Wolfram, "The wolfram functions site," <http://functions.wolfram.com>.
- [44] J. Zhang, W. Zeng, X. Li, Q. Sun, and K. P. Peppas, "New results on the fluctuating two-ray model with arbitrary fading parameters and its applications," *IEEE Trans. Veh. Technol.*, vol. 67, no. 3, pp. 2766–2770, Mar. 2017.
- [45] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [46] R. Girshick, "Fast r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448.
- [47] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Advances in neural information processing systems*, vol. 28, 2015.
- [48] D. Xu, Y. Zhu, C. B. Choy, and L. Fei-Fei, "Scene graph generation by iterative message passing," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5410–5419.
- [49] R. Zellers, M. Yatskar, S. Thomson, and Y. Choi, "Neural motifs: Scene graph parsing with global context," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 5831–5840.
- [50] H. Liu, N. Yan, M. Mortazavi, and B. Bhanu, "Fully convolutional scene graph generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 11 546–11 556.
- [51] A. Kroner, M. Senden, K. Driessens, and R. Goebel, "Contextual encoder–decoder network for visual saliency prediction," *Neural Netw.*, vol. 129, pp. 261–270, 2020.
- [52] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *International Conference on Learning Representations (ICLR)*, 2015.
- [53] X. Luo, Z. Liu, W. Wei, L. Ye, T. Zhang, L. Xu, and J. Wang, "Few-shot personalized saliency prediction using meta-learning," *Image and Vision Computing*, p. 104491, 2022.
- [54] W. Lu, X. Zhang, H. Lu, and F. Li, "Deep hierarchical encoding model for sentence semantic matching," *J. Vis. Commun. Image Represent.*, vol. 71, p. 102794, 2020.
- [55] F. Herrera, M. Lozano, and J. L. Verdegay, "Tackling real-coded genetic algorithms: Operators and tools for behavioural analysis," *Artif. Intell. Rev.*, vol. 12, no. 4, pp. 265–319, 1998.
- [56] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma *et al.*, "Visual genome: Connecting language and vision using crowdsourced dense image annotations," *International journal of computer vision*, vol. 123, no. 1, pp. 32–73, 2017.
- [57] M. Jiang, S. Huang, J. Duan, and Q. Zhao, "Salicon: Saliency in context," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, June 2015.

- [58] T. Judd, K. Ehinger, F. Durand, and A. Torralba, "Learning to predict where humans look," in *IEEE Int. Conf. Comput. Vis.*, 2009.
- [59] C. Yang, L. Zhang, H. Lu, X. Ruan, and M.-H. Yang, "Saliency detection via graph-based manifold ranking," in *IEEE Conf. Comput. Vis. Pattern Recognit.*, 2013, pp. 3166–3173.
- [60] P. Zhang, J. Zhang, K. P. Peppas, D. W. K. Ng, and B. Ai, "Dual-hop relaying communications over Fisher-snedecor  $\mathcal{F}$ -fading channels," *IEEE Trans. Commun.*, vol. 68, no. 5, pp. 2695–2710, May 2020.
- [61] H. Du, J. Zhang, K. P. Peppas, H. Zhao, B. Ai, and X. Zhang, "On the distribution of the ratio of products of Fisher-snedecor  $\mathcal{F}$  random variables and its applications," *IEEE Trans. Veh. Technol.*, vol. 69, no. 2, pp. 1855–1866, Feb. 2019.