

Generative AI for Semantic Communication: Architecture, Challenges, and Outlook

Le Xia, Yao Sun, Chengsi Liang, Lei Zhang, Muhammad Ali Imran, and Dusit Niyato

Abstract—Semantic communication (SemCom) is expected to be a core paradigm in future communication networks, yielding significant benefits in terms of spectrum resource saving and information interaction efficiency. However, the existing SemCom structure is limited by the lack of context-reasoning ability and background knowledge provisioning, which, therefore, motivates us to seek the potential of incorporating generative artificial intelligence (GAI) technologies with SemCom. Recognizing GAI's powerful capability in automating and creating valuable, diverse, and personalized multimodal content, this article first highlights the principal characteristics of the combination of GAI and SemCom along with their pertinent benefits and challenges. To tackle these challenges, we further propose a novel GAI-assisted SemCom network (GAI-SCN) framework in a cloud-edge-mobile design. Specifically, by employing global and local GAI models, our GAI-SCN enables multimodal semantic content provisioning, semantic-level joint-source-channel coding, and AIGC acquisition to maximize the efficiency and reliability of semantic reasoning and resource utilization. Afterward, we present a detailed implementation workflow of GAI-SCN, followed by corresponding initial simulations for performance evaluation in comparison with two benchmarks. Finally, we discuss several open issues and offer feasible solutions to unlock the full potential of GAI-SCN.

I. INTRODUCTION

Recently, semantic communication (SemCom) is popularized as an emerging paradigm that promises to significantly alleviate the scarcity of communication resources in future wireless networks [1]. This is mainly benefited from prosper advancement in deep learning (DL) technologies that can drive semantic encoding and decoding models to achieve efficient and high-quality semantic refinement on desired meaning with low spectrum consumption. Moreover, through equipping both ends of the transceiver with equivalent background knowledge, i.e., provisioning massive data samples to serve diverse artificial intelligence (AI) learning and prediction tasks [2], the implicit meaning in conveyed content can be recovered with ultra-low semantic errors even under harsh channel conditions. However, realizing such superiorities obviously poses a huge demand on data acquisition for constructing background knowledge and pre-training DL-driven semantic models. Meanwhile, considering that the achievable semantic performance is essentially confined by the quality of pre-training data used, existing SemCom systems still lack sufficient context reasoning capabilities, i.e., accurate semantic

calibration and recovery in transmitting multiple complex and coherent contextual fragments.

Fortunately, state-of-the-art (SOTA) generative AI (GAI) models, have lately emerged as killer applications in many verticals, promising to bring considerable productivity, innovation, and economic value to the real-world services as diverse as image synthesis, text generation, and drug discovery [3]. To be concrete, GAI leverages powerful DL algorithms, such as Transformer and diffusion to automate photorealistic and multimodal AI-generated content (AIGC) in response to user-provided prompts, while its fidelity and accuracy are contingent on adequate pre-training on billions of parameters in large language (e.g., ChatGPT) or image (e.g., Dall-E) models. Most importantly, GAI has immense abilities of context-reasoning and cross-modal content synthesis to generate high-quality and basically correct responses by successfully mimicking human's thinking and speaking patterns, enabling users feel that they are interacting with a real human-being rather than a dull machine [4]–[7]. Such a milestone innovation stimulates us to investigate the potential of applying AIGC into wireless SemCom, which hypothesizes to yield the below benefits.

Better SemCom Training Efficiency: Undoubtedly, the growing prosperity of SemCom is inseparable from colossal data resources for semantic model pre-training and background knowledge preparation targeting effective semantic interpretation [2]. To enable better SemCom training efficiency, GAI models are capable of producing vast multimodal content with a certain degree of authenticity and thus should be valuable materials to semantic training and background knowledge.

Enhanced Semantic Context-Reasoning: The historical AIGC can be stored online and retrieved easily to provision semantic coding models a better understanding for the context information, thereby offering significant context reasoning and semantic generalization. Furthermore, the creativity of GAI models can offer high-quality and precise content automatically for semantic interpretation in SemCom, thanks to the Prompt Engineer mechanism [8]. This ensures high semantic fidelity even if coding errors occur during joint-source-channel coding (JSCC) or physical signal transmission.

Higher Spectrum Utilization: Notice that most of AIGC can be produced via inputting only a few prompts, and the responses can be highly precise and specific if the prompts are well-crafted to align with a task-oriented communication. Hence, by utilizing well-trained GAI models, only several prompts need to be sent, instead of transmitting the whole source information, in each SemCom process to significantly deduct the required bandwidth resources while retaining the original meaning.

Le Xia, Yao Sun (*corresponding author*), Chengsi Liang, Lei Zhang, and Muhammad Ali Imran are with University of Glasgow, United Kingdom; Dusit Niyato is with Nanyang Technology University, Singapore.

Despite many ascendancies offered by the combination of GAI and SemCom, it still encounters several inevitable and thorny challenges that should be tackled before practical implementation. Among them, the paramount issue is how to deal with such considerable computing and storage resources required by these large GAI models. For instance, one of OpenAI's public large language models, called ChatGPT-3, comprises approximately 175 billion parameters in total [3], and hence it undoubtedly needs to consume colossal computing resources to operate the system. Another problem worth pointing out here lies in the reliability and latency aspects. AIGC is autonomously created that may lead to uncertainty to some extent, while introducing extra delay for data processing as well as data dissemination. To the best of our knowledge, only a few studies have explored the potential of incorporating GAI with SemCom. For example, the authors in [4] and [9] employed different deep GAI networks to enhance the perceptual quality and semantic reliability for SemCom. Similarly, [10] and [11] respectively sought the possibility of using GAI to quantify the semantic importance and adapt the end-to-end transmission rate. However, all of these works consider it from a transceiver design perspective in device-to-device scenarios, while relevant research for the system framework design in general cellular networks is still lacking. To this end, this article focuses on unleashing the full potential of GAI-assisted SemCom network (GAI-SCN) across the cloud-edge-mobile layers, and the main contributions are summarized in a nutshell as follows:

- We first present the basic concept of SemCom and four major types of GAI technologies, and then provide a comprehensive comparison among traditional communication, SemCom, and GAI-assisted SemCom. Next, we present the potential junctions between SemCom and GAI.
- We propose a novel GAI-SCN framework that integrates global and local GAI with semantic coding models in a collaborative cloud-edge-mobile design. Afterward, we showcase its viable implementation workflow consisting of three successive stages: Initial Network Preparation Stage, GAI-assisted SemCom Service Provisioning Stage, and Model Synchronization and Update Stage. Moreover, numerical results validate that the proposed framework can save a significant number of transmitted bits while maintaining high-precision semantic delivery compared with two benchmarks.
- Finally, several open issues with prospects of GAI-SCN are outlined, including device hardware limitations, inactive information sharing of users, and potential data tampering and privacy leakage.

II. WHEN SEMCOM MEETS GAI

In this section, we first introduce typical technologies of SemCom and GAI, as shown in Fig. 1, followed by several junctions between them identified and discussed in detail.

A. SemCom Systems

Compared with traditional communication of guaranteeing the precise reception of transmitted bits, the accurate delivery of semantics implied in desired messages becomes the

cornerstone of SemCom [11]. Taking an end-to-end SemCom system as the example, a transmitter first leverages background knowledge relevant to source messages to filter out irrelevant content and extract core features that only require fewer bits for transmission, the process of which is called semantic encoding. Once the receiver has the required knowledge, its local semantic interpreters are capable of accurately restoring the original meanings from the received bits, even with intolerable bit errors in data dissemination. This process is called semantic decoding. Consequently, efficient exchanges for the desired information with ultra-low semantic ambiguity can be achieved in SemCom under equivalent background knowledge, while significantly alleviating the resource scarcity problem.

B. Typical GAI Technologies

With proper pre-training and fine-tuning alongside extensive datasets, GAI excels in learning background knowledge and content structures from input training data, thereby generating outputs that closely resemble real-world samples [3]. In what follows, we briefly introduce four basic GAI technologies: generative adversarial networks (GANs), variational auto-encoders (VAEs), diffusion probabilistic models (DPMs), and flow-based generative models (FGMs). Note that other technologies like Transformer (related to ChatGPT) or the variants of these four (related to Dall-E) are equally essential as the indispensable components in the most of existing mainstream GAI models.

GANs: The GAN consists of two distinct neural networks: a generator and a discriminator. In the form of constant contestation, the goal of generator is to confuse the discriminator, while the discriminator should distinguish the samples generated by the generator from the real samples, until reaching a stable equilibrium [4]. Although the GAN enables GAI-SCNs to output a certain degree of accurate content, it is still not satisfactory enough in terms of semantic quality, generation diversity, and multimodal distribution problem learning.

VAEs: The VAE is a likelihood-based generative auto-encoder model, normally comprising of a multi-layer encoder and a symmetric decoder. By taking random training samples with a specific distribution as input, the encoder can regularize its coding distribution to ensure good properties of latent space, while the decoder maps from the latent space to the input space so as to produce new data points [5]. When it comes to the GAI-SCNs, efficient multi-task SemCom is envisioned to be realized with the support of VAEs.

DPMs: The DPM is a class of latent variable models, and its generation principle mainly includes two processes of forward diffusion and reverse diffusion. In the forward process, the input content is polluted in steps by introduced Gaussian noise. Afterward, the neural network in the reverse process should be trained to denoise the content blurred with noise, aiming to recover the original data. In comparison of GANs and VAEs, DPMs render remarkable superiority in semantic recovery tasks (e.g., image denoising, inpainting, and super-resolution), producing high-quality content and have better resistance to the risk of noise and interference [6].

FGMs: Differing from the previous models, FGMs are exact log-likelihood models with tractable sampling and latent

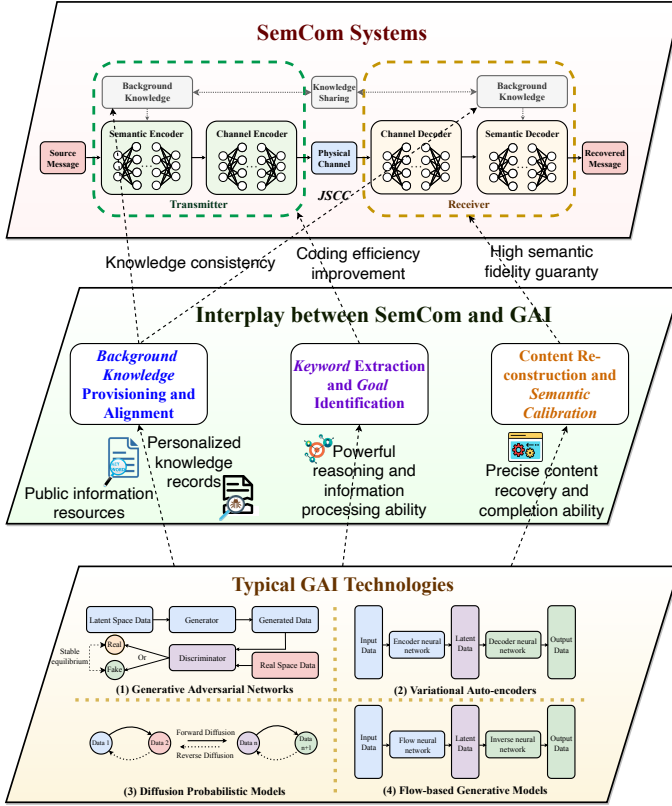


Fig. 1. Overview of SemCom systems and four types of typical GAI technologies along with three aspects of interplay between SemCom and GAI.

variable inference, which applies a bunch of reversible transformations to samples from the prior so that log-likelihoods of observations can be computed [7]. Besides, it leverages the change-of-variable law of probabilities to transform a simple distribution into a complex one, which greatly facilitates the semantic generation accuracy and communication efficiency.

Lesson Learned: As introduced, these SOTA GAI technologies are foreseeable to endow SemCom with a high level of semantic interpretation, reproducing and efficient content adjustment by quickly understanding input prompts in a human-thinking manner. However, the localization requirement for semantic models and the cloudification demand for large GAI models are inevitable, and thus it is meaningful to design a realistic GAI-SCN framework in order to yield more benefits from the communication perspective, such as lower consumption of resources, smarter semantic inference, and personalized SemCom services.

C. Interplay between SemCom and GAI

Based on the above introduction of SemCom systems and GAI technologies, herein we list three principal conjunctions between them with the corresponding elaboration.

Background Knowledge Provisioning and Alignment: Technically, GAI can be exploited as valuable data assets for users' background knowledge provisioning, which is roughly divided into two categories of global knowledge and personalized knowledge. For starters, global knowledge represents the common information publicly available to society (e.g.,

the content recorded in books, articles, videos, and other online sources), while personalized knowledge indicates users' personal information (e.g., language habits and communication style preferences). Thanks to sufficient pre-training, existing large GAI models like ChatGPT can easily and quickly retrieve global knowledge online to be stored as the common background knowledge of users. More importantly, such AI-generated global knowledge guarantees information consistency between any pair of communication parties, which ensures knowledge equivalence in SemCom. As for the personalized knowledge, it can be stipulated that GAI models store private conversations with users in the preparation stage of SemCom, by which their preferences can be analyzed in the background so as to offer personalized and customized AIGC according to local environments.

Keyword Extraction and Goal Identification: Recap that the core of SemCom is meaning delivery, for this purpose, GAI models are capable of extracting some keywords from the input long content, and the corresponding communication goal can be identified in a small amount of text (or the bounding box of objects in image) through excellent context-understanding ability. For instance, Bob uses a mobile device to share his personalized Paris-itinerary plan, detailing his week-long vacation schedule. In this case, several candidate keywords may appear like "One-week long vacation", "Bob's travel preferences", and "Landmarks and recommended food in Paris", etc., while the goal can be recognized as "Travel plan sharing from Bob". Accordingly, the global GAI model (e.g., a stable diffusion model) is trained specifically to restore the original content from received keywords and goals given background knowledge about Bob. Moreover, fewer communication resources (including wireless bandwidth and energy) are demanded, while the pressure from stringent latency requirements is relieved.

Content Reconstruction and Semantic Calibration: Co-operating with the keyword extraction function at the transmitter side, the GAI models deployed either in the core network or on the receiver side can realize content reconstruction according to different SemCom goals. Besides, a certain degree of semantic ambiguity may inevitably arise in the process of signal transmission and semantic interpretation, such as wording or sentence structure errors in the delivered text and blurred or partially missing tiles in the delivered images. To this end, the initially recovered content after semantic decoding can be input into some GAI models (like GPT-Neo) for fundamental and comprehensive semantic calibration to further improve the accuracy and reliability of SemCom.

Lesson Learned: Evidently, the GAI is promising to assist SemCom especially for task-oriented (i.e., meaning delivery-driven), high-capacity, and latency-insensitive services. Notably, although producing AIGC can consume some extra local processing time, the transmission delay is greatly reduced by SemCom in parallel, promising to significantly relieve the latency-cost pressure. Furthermore, we specially compare the system characteristics of GAI-assisted SemCom with traditional communication and SemCom alongside their respective benefits and limitations, as sketched in Fig. 2. Among them, traditional communication is built on bit-based source-and-

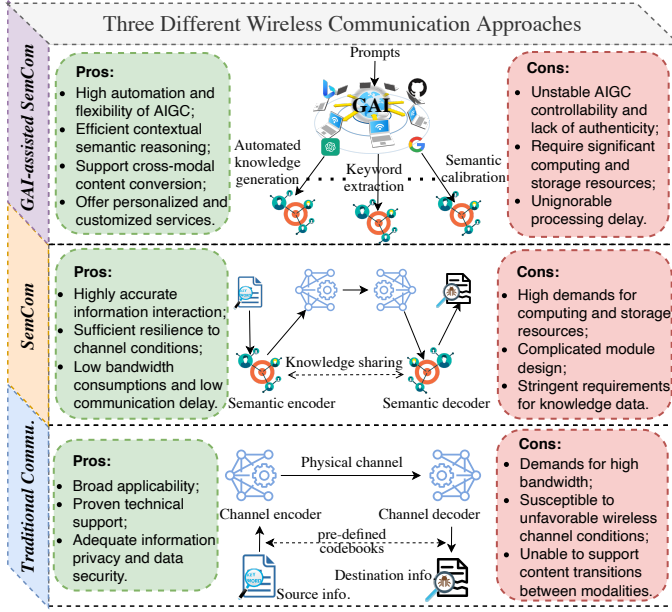


Fig. 2. Comparisons among three different approaches of GAI-SemCom, SemCom, and traditional communication in terms of their pros and cons.

channel coding following pre-defined and stringent codebooks, while SemCom integrates AI-based semantic encoder and decoder with equivalent background knowledge for accurate semantics-aware JSCC. On this basis, GAI-assisted SemCom takes full advantage of GAI technologies to significantly enhance semantic delivery efficiency and ease resource pressure.

III. GAI-ASSISTED SEMCOM NETWORK FRAMEWORK

In this section, we first propose a collaborative GAI-SCN framework, after which its implementation workflow is illustrated step by step.

A. Hierarchical Structure of GAI-SCN

Consider a SemCom-enabled cellular network scenario as demonstrated in Fig. 3, where there are multiple terminal devices (TDs) of senders and receivers within the coverage of base stations (BSs). Among them, multimodal SemCom services (e.g., text, image, and video) with specific communication goals consecutively arrive at each TD, and each BS acts as the service controller to efficiently schedule and coordinate the goal-oriented SemCom service provisioning. Additionally, a large GAI model (e.g., GPT-4 or Dall-E) is deployed in the cloud to complete computationally intensive tasks, while a small one (e.g., GPT-Neo) is embedded in the TD for coping with local lightweight service demands, such as customized content extraction and text generation services. Specially, we detail the proposed GAI-SCN from the perspective of a three-layer framework design.

Service Provisioning in the Mobile Layer: In the proposed GAI-SCN framework, each TD is equipped with a lightweight GAI model to realize context-aware keyword extraction and goal identification, taking full advantage of its reasoning ability in content understanding. For example, an open-source variant model of ChatGPT, called GPT-Neo, is a family of

Transformer-based language models developed by EleutherAI, which has a far smaller size than GPT-3 [12]. Through pre-training on large datasets (like Wikipedia and Common Crawl) and fine-tuning on user data (to be more personalized and customized), GPT-Neo is capable to be directly installed into TDs to precisely extract keywords of source information for users and identify the communication goal in only several words, by which fewer bits and smaller latency are consumed on the transmitter side for data transmission. As for the TDs on the receiver side, semantic decoders are deployed to recover the delivered meaning from obtained bits, after which the GPT-Neo can be further used for semantic calibration or language translation, etc.

JSCC Process in the Edge Layer: By enabling sufficient computing and storage ability in distributed edge servers, semantic encoders are considered to be deployed in the edge layer of GAI-SCN. On one hand, we exploit the JSCC method as it is able to greatly improve resilience and robustness of SemCom against various channel states, especially for the case with severe signal impairment, where the essential processes like pre-training, fine-tuning, and reasoning can be completely handled by edge computing servers. On the other hand, considering mass AIGC generated in the cloud or core network, it generally demands considerable computing resources for content pre-processing and semantic encoding before delivering them to users, which can be also offloaded to edge servers.

AIGC Acquisition in the Cloud Layer: Based on the above mobile and edge layer designs, the generation of AIGC becomes an indispensable process for smooth SemCom. In the proposed GAI-SCN, a centralized infrastructure, i.e., the remote cloud server, can support and run large GAI models like Google Bard or Microsoft Bing Chat. All preparation processes related to the AI model itself, such as pre-training and fine-tuning, are accomplished in the cloud. Likewise, the collection, analysis, and reasoning for multi-users' personal information and preferences can be realized in real time by leveraging the massive computing and storage resources of the cloud servers. Moreover, according to the specific keywords and goals uploaded from sender TDs, these large GAI models can quickly and correctly create the response content of desired modes. The rationale behind this is due to the historical conversation record, i.e., referred to as context, between users and GAI servers, which provides reference information for models' understanding and inference. Therefore, the original meaning implied in the delivered keywords can be recovered to intact text or image content.

Lesson Learned: By adopting the above collaborative GAI-SCN framework, the communication resource utilization and performance (e.g., semantic reasoning and generalization capabilities) can be maximized. Particularly, the cloud layer enables global GAI models to periodically coordinate all up-to-date interaction records uploaded from each user, whose personal communicating preferences can be predicted to provide semantic encoders in the edge layer with accurate and integral AIGC services. In parallel, the edge computing servers ensure high reliability and low latency for semantic processing and dissemination conditions via JSCC to provision a very high-

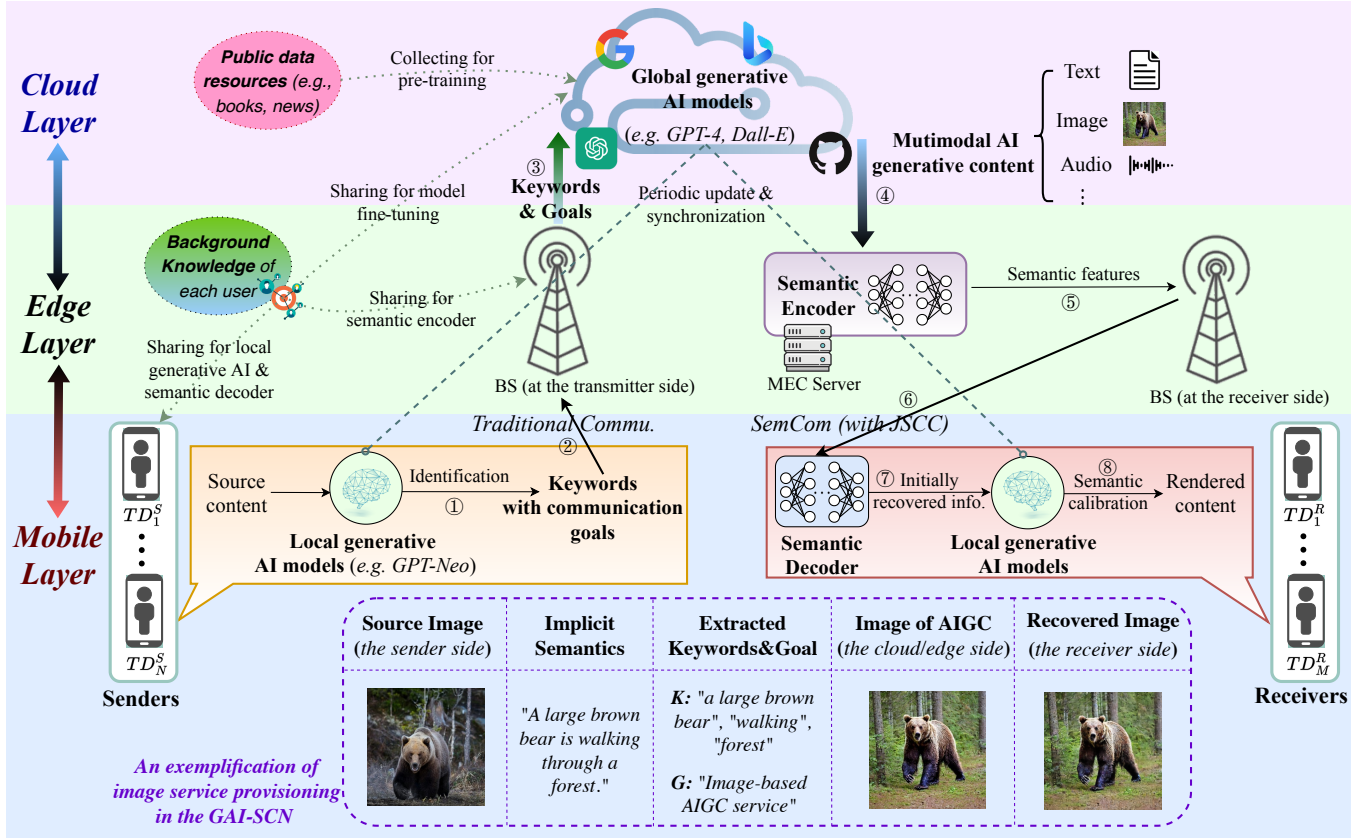


Fig. 3. Illustration of the proposed GAI-SCN framework in a collaborative cloud-edge-mobile design, where an exemplification of image service provisioning is presented.

quality user experience. Furthermore, light-weight GAI models deployed in the mobile layer can share users' personalized knowledge and keep it updated in real time. This guarantees consistency of AI's understanding for user preferences so as to prevent possible information mismatch between keyword extraction and recovered AIGC.

B. Implementation Workflow

After introducing the functionalities of the three layers in GAI-SCN, we revisit its key rationales and showcase the workflow below step-by-step to guide network designers to make proper changes on relevant protocols, as shown in Fig. 4.

Initial Network Preparation Stage: In this stage, the construction, pre-training and fine-tuning of GAI models are implemented collaboratively across the cloud and mobile layers. Note that all AI models are well pre-trained under society data public to users in the first place, like text from Wikipedia or books and images and videos from Instagram or YouTube, etc. Subsequently, considering the preference discrepancy between different users, personal data are collected by local GAI and then shared with global GAI in the cloud. Next, multiple parallel branch models corresponding to multiple users are created and fine-tuned based on personal data, where the branch model refers to a design pattern about multiple paths or subnetworks to make different predictions. Afterward, all parameters trained in each branch model are downloaded into the local light-weight GAI model related to each user, so as

to avoid the drawback of limited computing ability in mobile TDs. Apart from the above, the construction and joint pre-training of semantic encoders and decoders are also completed in advance given specific channel state information. Keeping in mind the requirement of knowledge equivalence in SemCom, knowledge sharing between encoders and decoders is required if the condition is triggered [2].

GAI-Assisted SemCom Service Provisioning Stage: Once all communication parties in the GAI-SCN are equipped with mature-trained AI models, the SemCom transmission services begin. The local GAI model first extracts keywords and identifies the explicit goal (referring to the exemplification in Fig. 3) for each user. Hereafter, such keywords and goals are uploaded via traditional communications to global GAI models in the cloud layer as input prompts, so that the corresponding branch AI model can create AIGC in line with the user preference to recover the original meaning with the original modality.

When it comes to the downlink side, each edge server detects the current wireless channel state and analyzes the received AIGC service-related signaling to select an appropriate pre-trained semantic encoder. As such, the entire AIGC is smoothly encoded into semantic features to be transmitted to the corresponding TD. Since semantic errors may still occur in the above JSCC process due to potential signal impairment and the limited computing power of TD, under the assistance of personalized knowledge and user preference, its local generative-AI model is utilized for further semantic

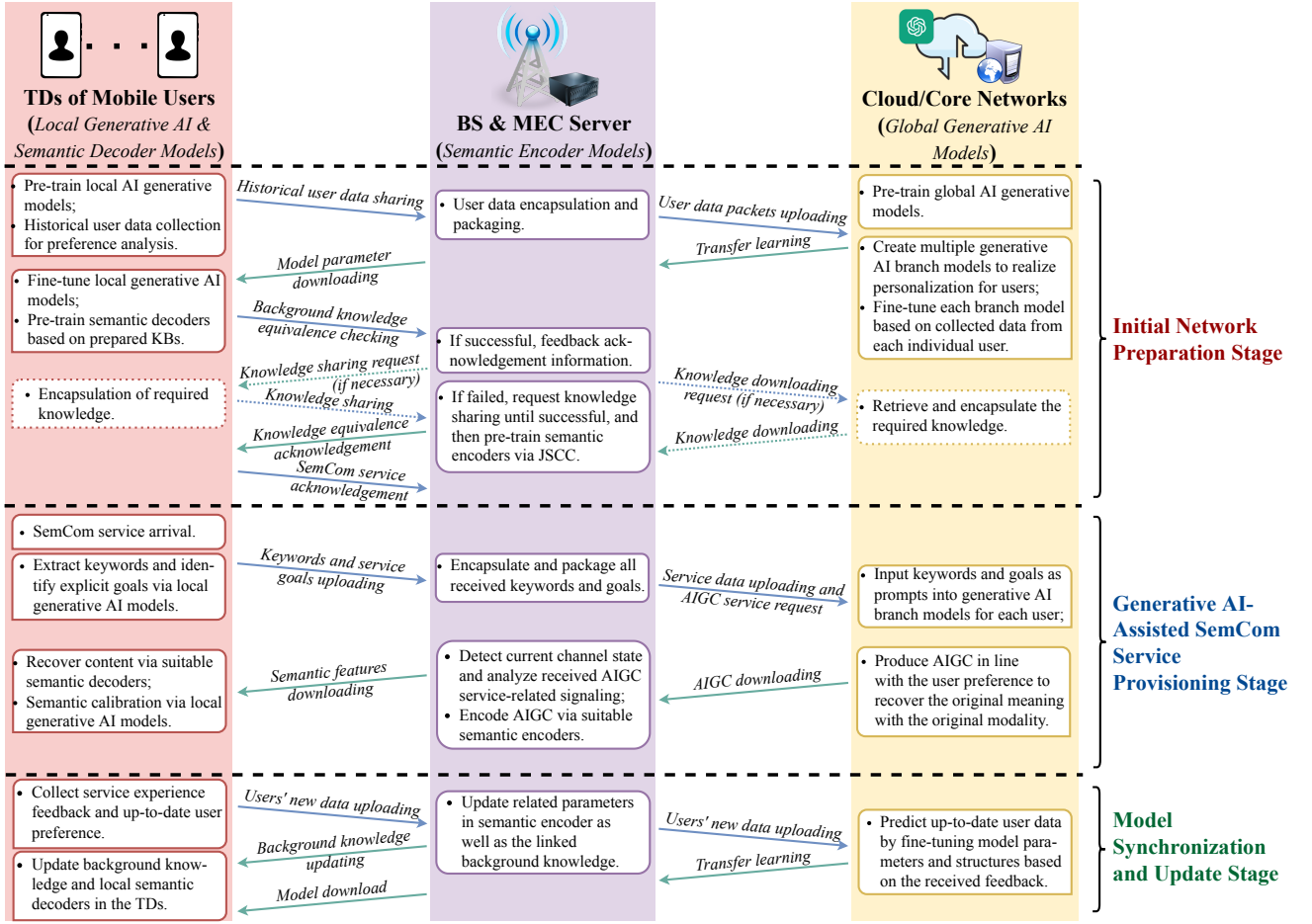


Fig. 4. A schematic diagram of implementing a complete round of semantic service provisioning in the GAI-SCN, including three successive stages of Initial Network Preparation, Generative AI-assisted SemCom Service Provisioning, and Model Synchronization and Update.

calibration, such as error correction for text content, color rendering for image content, and frame completion for video content, to enhance the resilience and robustness against semantic ambiguity for SemCom.

Model Synchronization and Update Stage: At the end of each GAI-assisted SemCom process, the service experience feedback (i.e., the user satisfaction regarding service performance, availability, and accessibility, etc.) is collected from each user and cached in the associated edge server. The feedback data as well as real-time user transmission data are synced periodically from the edge layer to the GAI models deployed in the other two layers. In this way, GAI analyzes and predicts users' up-to-date preferences by fine-tuning model parameters and structures. Meanwhile, the data are also fed back to semantic models to update the relative parameters as well as the linked background knowledge.

IV. CASE STUDY: IMAGE TRANSMISSION SERVICE PROVISIONING IN GAI-SCN

In this section, numerical results for a case study of image transmission are presented to evaluate the initial performance of the proposed GAI-SCN framework. For the simulation settings, an image-captioning model by combining ViT model with GPT-2 model [13] is exploited as the local GAI to

realize the image-to-text transformation as well as the keyword extraction and goal identification. Besides, we employ the latest text-to-image model called Stable Diffusion 2.1 [6] as the global GAI to create the AI-generated images from received prompts. As for the SemCom part, the main setups are proceeding as in the work [2], where an advanced deep convolutional network named Observation-Centric Sort and a Transformer-powered semantic decoder are leveraged for semantic segmentation and recovery, respectively. Meanwhile, all semantic models are trained based on the additive white Gaussian noise channel with a signal-to-noise ratio of 0 dB to transmit 327 images with different contents for testing. Finally, the Adam optimizer is adopted to train the neural networks in GAI-SCN with an initial learning rate of 5×10^{-4} based on the given image dataset. In parallel, for comparison purposes, we utilize two benchmarks: 1) A GAI-assisted traditional communication scheme, where the AIGC should be encoded into bits based on the prescribed coding rule [14] for precise image delivery; 2) A typical SemCom scheme [1], which transmits the original images solely via semantic coding models without any participation of GAI.

Figure 5 first demonstrates the performance of GAI-SCN by comparing the recovered images with the original ones under different numbers of observable objects that can be

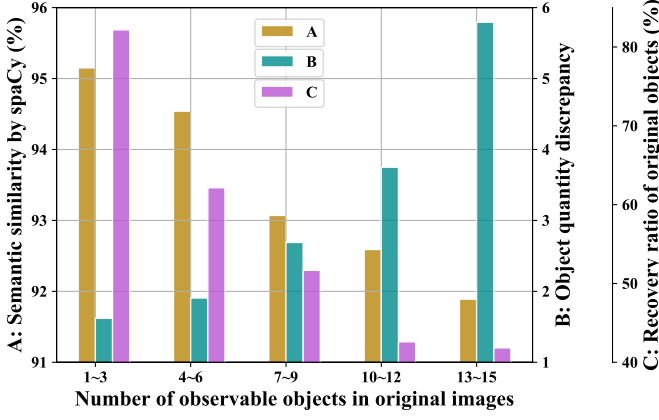


Fig. 5. Comparisons between original and recovered images by the proposed GAI-SCN framework in terms of three metrics: A) Semantic similarity by spaCy; B) Object quantity discrepancy; C) Recovery ratio of original objects.

TABLE I
AMOUNT OF BITS REQUIRED AND THE PSNR PERFORMANCE IN TRANSMITTING 300 IMAGES ON THE DOWNLINK (1024 * 1024 PIXELS)

Different image transmission schemes	Number of required bits for downlink	PSNR
GAI-assisted traditional communication [14]	1.28×10^5	28.05
SemCom [1]	5.99×10^4	28.25
GAI-SCN	3.03×10^4	28.64

detected and segmented (such as “bear”, “tree”, and “ground” as shown in the exemplification image in Fig. 3). Note that the appearances of the transmitted images may be not completely consistent with the received ones, however, the semantics implicit in the content should be our sole focus in measuring the system performance. Herein, we first test the semantic similarity performance measured by spaCy [15], where generally, the higher the spaCy score, the more accurate the recovered semantics. It is observed that the increasing number of objects in original images results in a decreasing semantic similarity, which is because that the higher complexity of images makes GAI more difficult to extract the keywords correctly as well as the image recovery. This phenomenon is also consistent with the recovery ratio performance of original objects, which metric shows the proportion of objects accurately recovered in the AIGC, and as the complexity of images rises, the average recovery ratio drops from 82.2% to 41.8% steadily. Moreover, we can see a lower object quantity discrepancy with fewer number of objects contained in the transmitted images. All of these trends above represent that semantic ambiguity is more likely to occur in regenerating more complex images due to confusing key object identification and a certain degree of semantic interference.

In addition, in order to validate the superiority of integrating SemCom with GAI, Table I shows the number of transmission bits required for three different schemes on the downlink side as well as the corresponding peak signal-to-noise ratio (PSNR) performance. Notably, all results in this table are based on the test of transmitting 300 images with a given size

of $1024 * 1024$ pixels. It can be found that the proposed GAI-SCN only requires 3.03×10^4 bits, which reduces 2.96×10^4 bits compared with the SemCom scheme and 9.77×10^4 bits with the GAI-assisted conventional scheme. Furthermore, the PSNR score obtained by our GAI-SCN maintains a very high level of 28.64, which is even slightly better than the other two approaches. This can be explained by the fact that the typical SemCom scheme starts the image delivery process from the uplink direction, which increases the risk of image sharpness loss, while the GAI-assisted traditional communication scheme is less resilient to harsh channel conditions compared with GAI-SCN. In summary, the above results demonstrate that our GAI-SCN can further save bandwidth resources while guaranteeing very high-quality SemCom service provisioning with accurate semantic delivery.

V. OPEN RESEARCH ISSUES AND OUTLOOKS

In this section, we list several thorny issues and outlooks that can be highlighted as future research directions in the GAI-SCN.

Limited Device Resources for Supporting AI Modules:

In the GAI-SCN, sophisticated AI-enabled computing modules (including local GAI models and SemCom coding models) need distributed implementation at each TD, imposing a heavy burden on its inherently limited device resources (like storage, memory, computational units, and battery power). To make it practically implementable, advanced model compression and acceleration technologies, such as knowledge distillation, parameter pruning and quantization, are promising to efficiently drop the complexity and size of AI networks with an affordable cost of performance degradation.

Randomness of Content Rendered in GAI-SCN: Notice that the appearance of AIGC output from the cloud GAI may vary even given the same keywords and goals. Besides, the representation of semantics recovered via the semantic decoder also have uncertainty to some extent, due to the knowledge mismatching or semantic errors. Therefore, granularity tuning on keyword extraction and subsequent semantic calibration deserve further investigation to tackle such randomness.

Inactive Sharing of Background Knowledge and Personal Preferences: Since the prerequisite of customized AIGC and SemCom services mainly lies in the proactive sharing of users’ personal preferences and background knowledge, devising a scores- or rewards-based incentive mechanism, such as delegated proof of stake-based blockchain, is necessary to attract users to spontaneously contribute personal data to the upgrading of GAI-SCN, where potential reward alternatives include social welfare and tech benefits, etc.

VI. CONCLUSIONS

This article explored the potential of applying AIGC into SemCom for service provisioning, where we first showcased the development of SemCom and GAI technologies with their integration cases, and then proposed the GAI-SCN framework. Specially, the collaborative cloud-edge-mobile structure was well-devised to incorporate both global and local GAI models with the JSCC process, which not only enables efficient and

high-quality meaning delivery, but also significantly reduces transmission traffic as well as latency. Moreover, implementation alongside initial simulations was provided, followed by associated open issues and corresponding solutions. We hope that our GAI-SCN serves as a pioneer in facilitating communication resource usage as well as user experience for futuristic context-aware and GAI-based wireless SemCom networks.

REFERENCES

- [1] D. Huang, X. Tao, F. Gao, and J. Lu, "Deep Learning-Based Image Semantic Coding for Semantic Communications," in *2021 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2021, pp. 1–6.
- [2] L. Xia, Y. Sun, C. Liang, D. Feng, R. Cheng, Y. Yang, and M. A. Imran, "WiserVR: Semantic Communication Enabled Wireless Virtual Reality Delivery," *IEEE Wireless Communications*, vol. 30, no. 2, pp. 32–39, 2023.
- [3] M. Xu, H. Du, D. Niyato, J. Kang, Z. Xiong, S. Mao, Z. Han, A. Jamalipour, D. I. Kim, V. Leung *et al.*, "Unleashing the Power of Edge-Cloud Generative AI in Mobile Networks: A Survey of AIGC Services," *IEEE Communications Surveys & Tutorials*, 2024.
- [4] E. Erdemir, T.-Y. Tung, P. L. Dragotti, and D. Gündüz, "Generative Joint Source-Channel Coding for Semantic Image Transmission," *IEEE Journal on Selected Areas in Communications*, 2023.
- [5] D. P. Kingma, M. Welling *et al.*, "An Introduction to Variational Autoencoders," *Foundations and Trends® in Machine Learning*, vol. 12, no. 4, pp. 307–392, 2019.
- [6] Z. J. Wang, E. Montoya, D. Munechika, H. Yang, B. Hoover, and D. H. Chau, "DiffusionDB: A Large-scale Prompt Gallery Dataset for Text-to-Image Generative Models," *arXiv preprint arXiv:2210.14896*, 2022.
- [7] D. Rezende and S. Mohamed, "Variational Inference with Normalizing Flows," in *International Conference on Machine Learning*. PMLR, 2015, pp. 1530–1538.
- [8] V. Liu and L. B. Chilton, "Design Guidelines for Prompt Engineering Text-to-Image Generative Models," in *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, 2022, pp. 1–23.
- [9] C. K. Thomas and W. Saad, "Neuro-Symbolic Causal Reasoning Meets Signaling Game for Emergent Semantic Communications," *IEEE Transactions on Wireless Communications*, 2023.
- [10] S. Guo, Y. Wang, S. Li, and N. Saeed, "Semantic Importance-Aware Communications Using Pre-Trained Language Models," *IEEE Communications Letters*, 2023.
- [11] S. Barbarossa, D. Communiello, E. Grassucci, F. Pezone, S. Sardellitti, and P. Di Lorenzo, "Semantic Communications Based on Adaptive Generative Models and Information Bottleneck," *IEEE Communications Magazine*, vol. 61, no. 11, pp. 36–41, 2023.
- [12] R. Kashyap, V. Kashyap *et al.*, "GPT-Neo for Commonsense Reasoning-A Theoretical and Practical Lens," *arXiv preprint arXiv:2211.15593*, 2022.
- [13] Y. Shen, K. Song, X. Tan, D. Li, W. Lu, and Y. Zhuang, "HuggingGPT: Solving AI Tasks with ChatGPT and its Friends in HuggingFace," *arXiv preprint arXiv:2303.17580*, 2023.
- [14] Z. Cai, J. Hao, P. Tan, S. Sun, and P. Chin, "Efficient Encoding of IEEE 802.11 n LDPC Codes," *Electronics Letters*, vol. 42, no. 25, p. 1, 2006.
- [15] Y. Vasiliev, *Natural Language Processing with Python and spaCy: A Practical Introduction..* No Starch Press, 2020.

Le Xia (l.xia.2@research.gla.ac.uk) is currently pursuing his Ph.D degree with the James Watt School of Engineering, University of Glasgow, UK. His research interests include semantic communications, intelligent vehicular networks, and resource management in next-generation wireless networks.

Yao Sun (Yao.Sun@glasgow.ac.uk) is currently a Lecturer with the James Watt School of Engineering, the University of Glasgow, UK. His research interests include semantic communications, intelligent wireless networking, and wireless blockchain system.

Chengsi Liang (23578751@student.gla.ac.uk) is currently pursuing her Ph.D. degree with the James Watt School of Engineering, University of Glasgow, UK. Her research interest includes semantic communication and networking.

Lei Zhang (Lei.Zhang@glasgow.ac.uk) is a Professor at the University of Glasgow. He has academia and industry combined research experience on wireless communications and networks, and distributed systems for IoT, blockchain, autonomous systems. He is the founding Chair of IEEE Special Interest Group on Wireless Blockchain Networks in Cognitive Networks Technical Committee.

Muhammad Ali Imran (Muhammad.Imran@glasgow.ac.uk) is a Professor of communication systems with the University of Glasgow, UK, and a Dean with Glasgow College UESTC. He is also an Affiliate Professor with the University of Oklahoma, USA, and a Visiting Professor at University of Surrey, UK. He has over 20 years of combined academic and industry experience with several leading roles in multi-million pounds funded projects.

Dusit Niyato (dniyato@ntu.edu.sg) is a professor in the School of Computer Science and Engineering, at Nanyang Technological University, Singapore. He received Ph.D. in Electrical and Computer Engineering from the University of Manitoba, Canada in 2008. He has published more than 400 technical papers in the areas of wireless and mobile computing, sustainability, edge intelligence, decentralized machine learning, and incentive mechanism design. He was a Distinguished Lecturer of the IEEE Communications Society from 2016 to 2017. He is a Fellow of IEEE.