

Semantic Information Extraction for Text Data with Probability Graph

Zhouxiang Zhao*, Zhaohui Yang*^{†‡}, Ye Hu[§], Licheng Lin*, and Zhaoyang Zhang*[‡]

*College of Information Science and Electronic Engineering, Zhejiang University, Hangzhou, China

[†] Zhejiang Lab, Hangzhou, China

[‡]Zhejiang Provincial Key Laboratory of Info. Proc., Commun. & Netw. (IPCAN), Hangzhou, China

[§]Department of Industrial and System Engineering, University of Miami, Coral Gables, USA

E-mails: zhouxiangzhao@zju.edu.cn, yang_zhaohui@zju.edu.cn, yxh1096@miami.edu,

linlicheng@zju.edu.cn, ning_ming@zju.edu.cn

Abstract—In this paper, the problem of semantic information extraction for resource constrained text data transmission is studied. In the considered model, a sequence of text data need to be transmitted within a communication resource-constrained network, which only allows limited data transmission. Thus, at the transmitter, the original text data is extracted with natural language processing techniques. Then, the extracted semantic information is captured in a knowledge graph. An additional probability dimension is introduced in this graph to capture the importance of each information. This semantic information extraction problem is posed as an optimization framework whose goal is to extract most important semantic information for transmission. To find an optimal solution for this problem, a Floyd’s algorithm based solution coupled with an efficient sorting mechanism is proposed. Numerical results testify the effectiveness of the proposed algorithm with regards to two novel performance metrics including semantic uncertainty and semantic similarity.

Index Terms—Semantic information extraction, knowledge graph, probability graph, semantic communication.

I. INTRODUCTION

Over the past few decades, the development of mobile communication technology has greatly contributed to the progress of human society. In the 1940s, Shannon proposed information theory [1], which focused on quantifying the maximum data transmission rate that a communication channel could support. Guided by this fundamental theory, most existing communication systems have been designed based on metrics which concentrate on transmission rate. With the rapid increase in demand for intelligent applications of wireless communication, the future communication network will change from a traditional architecture that simply pursues high transmission rate to a new architecture that is oriented to complete tasks efficiently. There will be more and more tasks require low latency and high efficiency in future mobile information networks, which poses a huge challenge [2], [3] to the existing mobile communication systems. This trend gives rise to a new communication paradigm, called *semantic communication*. Shannon’s theory did not give importance to the semantic information of the data. It was considered to be largely irrelevant to communication. The concept of semantic

communication was introduced by Shannon’s collaborator Warren Weaver. Semantic communication is a new architecture that can integrate user needs and information meaning into the communication process. Traditional communication systems are primarily dedicated to reliably transmitting bit streams. It does not know the meaning of the message or what the goal of the message exchange is. Semantic communication systems, on the other hand, are dedicated to conveying the meaning in the message [4], [5] and can significantly reduce the required transmission channel bandwidth [6]. Semantic communication is expected to be a key technology for the future 6G mobile communication systems [7], [8].

Recently, several works studied a number of problems related to semantic communications. The authors in [9] proposed a deep learning based semantic communication system for text transmission. To measure the performance of semantic communication, the authors also designed a new metric called sentence similarity. In [10], a deep learning-enabled semantic communication system for speech signals was designed. In order to improve the recovery accuracy of speech signals, especially for the essential information, it was developed based on an attention mechanism by utilizing a squeeze-and-excitation network. The work in [11] enabled the transmission of audio semantic information which captures the contextual features of audio signals. To extract the semantic information from audio signals, a wave to vector architecture based autoencoder that consists of convolutional neural networks was proposed. The authors in [12] proposed a lite distributed semantic communication system based on deep learning, named L-DeepSC, for text transmission with low complexity, where the data transmission from the IoT devices to the cloud/edge works at the semantic level to improve transmission efficiency. However, these existing works in [9]–[12] did not consider a flexible design that has the ability to adapt the content of the transmission in accordance with the communication environment.

The main contribution of this paper is a novel probability graph based important information extraction method, which enables the selection of the most important semantic information with in resource constrained networks. The key contributions are listed as follows:

- We propose a novel framework to extract important

This work is supported by Zhejiang Lab Program under grant K2023QA0AL02, Zhejiang Science and Technology Program under grant 2023C01021, National Natural Science Foundation of China under Grants U20A20158, and 61725104.

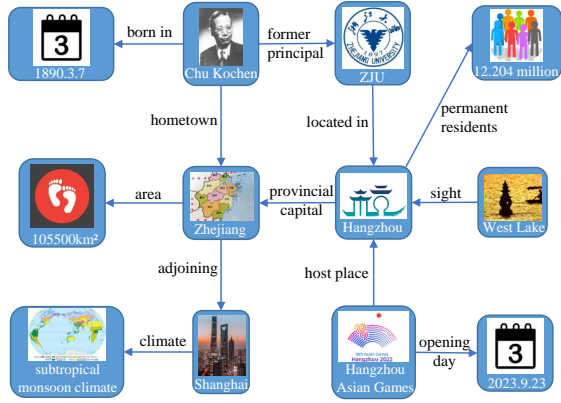


Fig. 1. An example of knowledge graph.

information for semantic text data transmission under communication resource constrains. In particular, we capture the semantic information in a knowledge graph, and introduce an additional relation probability dimension in the graph to capture the importance of the information. We formulate this information extraction problem in an optimization framework, and seek to find the optimal strategy that finds the most important information to extract and transmit.

- To guarantee the effectiveness of the proposed semantic information extraction solution, we introduce two evaluation metrics: semantic uncertainty, which measure the clarity of the semantic triples, and semantic similarity, which measures the similarity between two texts. Numerical results show that the proposed algorithm's effectiveness in regards to these two metrics.

Semantic communication usually relies on semantic information extraction, which can remove redundant information in the original data that is less relevant to the task. There are many methods to extract semantic information, and one of them is knowledge graph. An example of knowledge graph is shown in Fig. 1. Knowledge graph can carry a large amount of information with a small amount of data [13]. Therefore, it is a frequently used tool for semantic information extraction [14].

To our knowledge, few research has combined knowledge graph and random theory with wireless communication scenarios. For example, in a wireless resource-constrained scenario, after obtaining a knowledge graph from text data, if there are no sufficient communication resources to transmit the complete knowledge graph within a limited time delay, it is necessary to compress the original knowledge graph. However, there are few such knowledge graph compression methods in the existing research.

We adapt the work in [15] to generate semantic triples from text data and the confidences of each triple. Then, we propose an algorithm based on probability graph to select a part of the raw triples, which realizes the compression of knowledge graph.

II. SYSTEM MODEL

In order to achieve semantic information extraction of text data and reduce the burden of communication networks, this paper uses information extraction techniques in natural language processing (NLP) to extract semantic triples from text data. Secondly, to measure the performance of our algorithm, we propose two metrics called semantic uncertainty and semantic similarity. Those parts will be elaborated in the following separately.

A. Semantic Information Extraction

This step will first perform named entity recognition on the text data and output the tagged text. Then, the tagged text is fed into the relation extraction model to extract the relations among entities. The final extracted information is represented by a triple (*head, relation, tail*). In performing the relation extraction task, a relation set consisting of several specific relations is available in advance, and the model finally outputs the probability of each specific relation among the entities. Currently, relation extraction still faces some challenges with an accuracy of about 80% [16]. Therefore, the initial triples obtained from text data cannot all be correct, and subsequent steps are needed to minimize the impact of inaccurate triples.

In the studied model, we use w_n to represent a word, a symbol, or a punctuation in the text data. Hereinafter, w_n is called a token. Based on this, the original text data can be represented by an ordered sequence of tokens as shown below:

$$L = \{w_1, w_2, \dots, w_n, \dots, w_N\}, \quad \forall w_n \in \mathcal{V}, \quad (1)$$

where \mathcal{V} is the vocabulary and N is the number of tokens in L . For instance, assuming that the text data need to be transmitted is "Apple is a kind of fruit.". Hence, in this example, $L = \{[Apple], [is], [a], [kind], [of], [fruit], [.] \}$, where $w_1 = [Apple]$, $w_2 = [is]$, $w_3 = [a]$, $w_4 = [kind]$, $w_5 = [of]$, $w_6 = [fruit]$, $w_7 = [.]$.

In the applied model, the semantic information extracted from text data is represented by knowledge graph. A knowledge graph is a structured representation of facts, which consists of entities and relations. The knowledge in a knowledge graph can be represented by a triple (*head, relation, tail*). The knowledge graph consists of a series of nodes and edges [17]. To be more specific, the node in knowledge graph is called an *entity*, representing the object or concept in the real world. Define entity j in text data L as e_j , which consists of a series of tokens in L . The edge in knowledge graph represents the *relation* of two entities. Define the relation between entity pair (e_j, e_k) as $r_{jk} \in \mathcal{R}$, where \mathcal{R} is the relation set.

According to existing NLP techniques, the probability that all the relations in the relation set corresponding to each triple can be obtained. Further, the corresponding entropy, denoted as h_{jk} , can be obtained using the following formula:

$$h_{jk} = - \sum_{i=1}^A (p_{jk_i} \log_2 p_{jk_i}), \quad (2)$$

where A is the total number of relations in the relation set \mathcal{R} , and p_{jk_i} is the confidence of each relation for entity pair (e_j, e_k) .

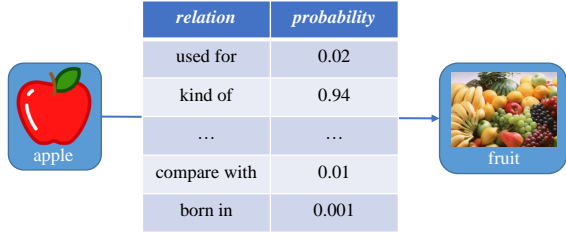


Fig. 2. Each edge between two entities consists of several relation probabilities.

The smaller the h is, the more explicit the relations are, and the larger the h is, the vice versa. Therefore, it may be appropriate to expand and write the triples in the form of $(e_j, r_{jk}, e_k, h_{jk})$.

Because the dimension of entropy is added based on probability, the result obtained in this step is not a traditional triple knowledge graph, but a probability graph version of the knowledge graph. The relation between each two nodes is not absolute, but consists of several relations corresponding to probabilities jointly, as shown in Fig. 2.

According to the knowledge graph, the semantic information in the text data L can be written in the following form:

$$\mathbb{G} = \{\varepsilon^1, \varepsilon^2, \dots, \varepsilon^g, \dots, \varepsilon^G\}, \quad (3)$$

where $\varepsilon^g = (e_j^g, r_{jk}^g, e_k^g, h_{jk}^g)$, $j \neq k$, and G is the total number of quadruples in \mathbb{G} . After obtaining the above knowledge graph, we use our algorithm to compress it and extract the relatively important information from it. The extracted important semantic information is represented as:

$$\mathbb{G}' = \{\varepsilon'^1, \varepsilon'^2, \dots, \varepsilon'^h, \dots, \varepsilon'^H\}, \quad (4)$$

where H is the total number of triples after compressing.

B. Semantic Uncertainty

To evaluate the quality of the semantic information extraction, we proposed a metric called *semantic uncertainty*. Based on the probability, the semantic uncertainty can reflect the ambiguity of the semantic information. We define the semantic uncertainty as:

$$SU = \sum_{h \in \mathbb{G}'} h, \quad (5)$$

where h is the entropy of each selected quadruple.

As we all know, entropy can measure the uncertainty of an event, the greater the uncertainty, the greater the entropy. Here, h represents the uncertainty of a triple in the knowledge graph, the greater the h , the less certain the semantic relation. Therefore, the lower the SU , the lower the uncertainty of the semantic information and the better the extraction quality.

C. Semantic Similarity

Another way to evaluate the quality of semantic extraction is to recover the extracted triples into text and then compare this recovered text with the original text to evaluate the similarity between them. The text recovery can be done with the help of Generative Pre-trained Transformer [18].

Thus, we proposed *semantic similarity* to measure the semantic similarity between two texts. The semantic similarity consists of two parts: semantic accuracy and semantic completeness. The semantic accuracy is defined as follows:

$$A(\mathbb{G}') = \frac{\sum_{m=1}^M \min\{\sigma(L'(\mathbb{G}'), e_m), \sigma(L, e_m)\}}{\sum_{m=1}^M \sigma(L'(\mathbb{G}'), e_m)}, \quad (6)$$

where M is the total number of different entities in the extracted triples, $\sigma(L'(\mathbb{G}'), e_m)$ is the number of occurrences of the entity e_m in the recovered text $L'(\mathbb{G}')$, and $\sigma(L, e_m)$ is the number of occurrences of the entity e_m in the original text L . The semantic completeness is defined as follows:

$$C(\mathbb{G}') = \frac{\sum_{m=1}^M \min\{\sigma(L'(\mathbb{G}'), e_m), \sigma(L, e_m)\}}{\sum_{m=1}^M \sigma(L, e_m)}. \quad (7)$$

Based on equation (6) and (7), the semantic similarity can be written as:

$$SS(\mathbb{G}') = \theta(\mathbb{G}') \frac{A(\mathbb{G}')C(\mathbb{G}')}{\varphi A(\mathbb{G}') + (1 - \varphi)C(\mathbb{G}')}, \quad (8)$$

where φ is a parameter that regulates the contribution of semantic accuracy and semantic completeness to semantic similarity, $\varphi \in [0, 1]$. A larger φ will increase the impact of semantic accuracy on semantic similarity, and vice versa. $\theta(\mathbb{G}')$ is defined as:

$$\theta(\mathbb{G}') = \sum_{p \in \mathbb{G}'} p, \quad (9)$$

where p is the largest relation probability in each extracted triples.

Specifically, for a certain entity e_m corresponding to the recovered text $L'(\mathbb{G}')$, if it appears more times in the recovered text $L'(\mathbb{G}')$ than in the original text L , then the semantic accuracy decreases; conversely, the semantic completeness decreases.

III. KNOWLEDGE GRAPH IMPORTANT INFORMATION EXTRACTION ALGORITHM

To further extract the backbone information with high confidence, this paper uses techniques which are related to probability graph to select the best combination.

A piece of text data usually has a central concept around which the rest of the text basically revolves. In knowledge graph, we call the central concept of the corresponding text as initial node. The scheme to determine initial node e_1 is as follows: the node with the most occurrences in \mathbb{G} is selected as the initial node; if it is not unique, we choose the node that appears earlier in the text.

Next, the relational distance $d(e_i, e_j)$ is defined as the minimum number of edges to be traversed from e_i to e_j , when the directionality of edges is not considered. The relational distance can measure the degree of association between two entities.

For example, in Fig. 3, there are five edges connected at the node *Hangzhou*, which means node *Hangzhou* appears five times in \mathbb{G} . Observing the other nodes, we find that none of them has more occurrences than node *Hangzhou*, so

head	relation	tail	entropy
Bruce Lee	born in	1940/11/27	0.0092
Bruce Lee	birthplace	San Francisco	0.1836
Bruce Lee	star in	Enter the Dragon	1.3390
Bruce Lee	star in	Dragon Crossing	1.2914
Bruce Lee	wife	Linda Emery	1.4416
Bruce Lee	nationality	USA	0.4132
Bruce Lee	graduation institution	University of Washington	1.2480
Linda Emery	nationality	USA	0.4812
Linda Emery	graduation institution	University of Washington	1.8124
USA	graduation institution	University of Washington	1.0055
Bruce Lee	father of	Brandon Lee	1.6349
Brandon Lee	born in	1965/2/1	0.0239
Bruce Lee	wife	Shannon Lee	1.5973
Shannon Lee	born in	1969/4/19	0.0102
Bruce Lee	author	Jeet Kune Do	2.3385
Hoi-Chuen Lee	born in	1901/2/4	0.0079
Hoi-Chuen Lee	father of	Bruce Lee	1.3843
Hoi-Chuen Lee	birthplace	Shunde	0.0799
Hoi-Chuen Lee	nationality	China	0.0184
Brandon Lee	star in	Spirit Of The Dragon	1.1067
Brandon Lee	star in	The Crow	0.8614
Shannon Lee	star in	Dragon: The Bruce Lee Story	0.1709
Shannon Lee	star in	Enter the Eagles	0.1667
Shannon Lee	producer	The Legend of Bruce Lee	1.6201

Fig. 5. Original quadruples obtained from the text data about “Bruce Lee”, and the selected quadruples using our algorithm when $K = 0.5$ and $D = 2$. Note that the selected quadruples are marked in green.

simulation, and the result is shown in Fig. 5. As we can see, the selected quadruples have low entropy value, which means they have relatively more explicit semantic relations according to the text.

As mentioned above, the compression coefficient K determines the compression ratio of the original knowledge graph. In order to find out the relationship between the compression coefficient K and the semantic uncertainty SU , we changed K from 0.1 to 1 and calculated SU respectively. For comparison, we considered several baselines: 1) chooses quadruples randomly, 2) extracts by the number of occurrences of the entity from most to least, 3) extracts by the number of occurrences of the entity from least to most, 4) extracts the triples in order of appearance from front to back, and 5) extracts the triples in order of appearance from back to front. The result is shown in Fig. 6.

In Fig. 6, we can see that, in our proposed algorithm, as the compression coefficient K increases, the semantic uncertainty

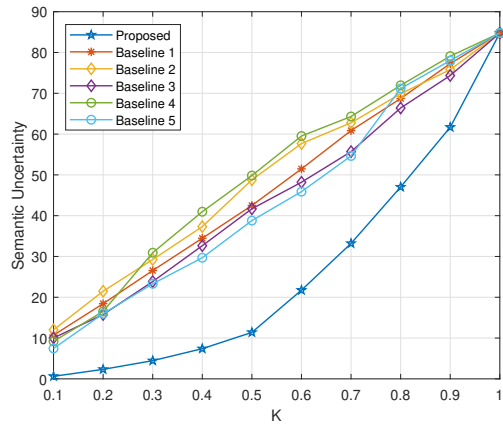


Fig. 6. The relationship between the compression coefficient K and the semantic uncertainty SU when the maximum depth $D = 2$. The result of baseline 1 is the average of 100 runs.

SU increases. This is because as the compression coefficient K increases, the total number of selected quadruples increases, and every quadruple will bring additional semantic uncertainty to SU . Fig. 6 also shows that the proposed algorithm always has lower SU than baseline algorithms, except when $K = 1$. In the case of $K = 1$, all original quadruples are chosen, so the SU value of all algorithms are the same. The result shows that our proposed algorithm can extract more explicit information from the original knowledge graph.

Original text	Original text		
Chosen triples $K = 0.5$ $D = 2$	Bruce Lee	born in	1940/11/27
	Bruce Lee	birthplace	San Francisco
	Bruce Lee	nationality	USA
	Linda Emery	nationality	USA
	Brandon Lee	born in	1965/2/1
	Shannon Lee	born in	1969/4/19
	Hoi-Chuen Lee	born in	1901/2/4
	Hoi-Chuen Lee	birthplace	Shunde
	Hoi-Chuen Lee	nationality	China
	Brandon Lee	star in	The Crow
	Shannon Lee	star in	Dragon: The Bruce Lee Story
	Shannon Lee	star in	Enter the Eagles
Recovered text	<p>Bruce Lee was born on November 27, 1940, in San Francisco, USA. He was a USA national. His wife, Linda Emery, is also a USA national. They have a son named Brandon Lee, who was born on February 1, 1965. Bruce Lee's father, Lee Hoi-chuen, was born on February 4, 1901, in Shunde, Guangdong, China. He was from China. Brandon Lee starred in the movie "The Crow". Shannon Lee was born on April 19, 1969. She starred in movies like "Dragon: The Bruce Lee Story" and "Enter the Eagles".</p>		

Fig. 7. An example of the original text, extracted triples, and corresponding recovered text. The entities in chosen triples are marked in red.

Further, the quality of semantic extraction is evaluated by semantic similarity after recovering the extracted triples into text. An example of the original text, extracted triples, and corresponding recovered text is presented in Fig. 7. The result in Fig. 8 demonstrates that an increase in the compression coefficient K leads to an increase in semantic similarity SS . This is because the selection of triples for text recovery increases with higher compression coefficients, resulting in a greater similarity between the recovered text and the original text. The front segment of the proposed algorithm's semantic similarity increases more rapidly than the end segment because the triples used in the former are of higher semantic quality and provide more accurate semantic information. Although the gap between the proposed algorithm and other algorithms decreases after the triples are recovered into text, the proposed algorithm still retains some advantages.

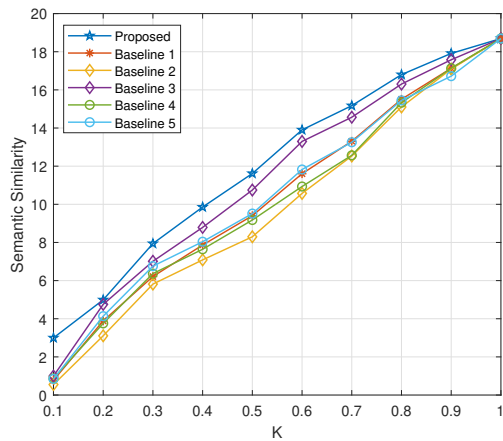


Fig. 8. The relationship between the compression coefficient K and the semantic similarity SS when the maximum depth $D = 2$.

V. CONCLUSION

In this paper, we leveraged knowledge graphs to model the semantic information of textual data and presented a novel algorithm for extracting important information from these graphs. To measure the performance of our approach, we introduced two metrics: semantic uncertainty, which reflects the clarity of the semantic triples, and semantic similarity, which measures the similarity between the recovered text and the original text. Our approach introduced the concept of relational probability into knowledge graphs and extends traditional triples into quadruples. We also proposed two parameters, the compression coefficient K and the maximum depth D , which determine the number of selected quadruples and the compactness of the extracted information, respectively. We then formulated an optimization problem that aims to minimize the information entropy of the selected semantic information while satisfying the constraints of the compression coefficient and maximum depth. Our simulation results demonstrate the effectiveness of the proposed knowledge graph important information extraction algorithm.

With the ability to effectively reduce the amount of data that needs to be transmitted while preserving the optimal textual

semantics, our algorithm can play an important role when the communication resources are severely limited. This capability becomes especially valuable in satellite communications, underwater sensor networks and other contexts with constrained communication resources.

Looking ahead, an intriguing and challenging research direction is the joint communication and computation resource allocation while incorporating probability graph of semantic information extraction.

REFERENCES

- [1] C. E. Shannon, "A mathematical theory of communication," *The Bell system technical journal*, vol. 27, no. 3, pp. 379–423, 1948.
- [2] D. Gündüz, Z. Qin, I. E. Aguerri, H. S. Dhillon, Z. Yang, A. Yener, K. K. Wong, and C.-B. Chae, "Beyond transmitting bits: Context, semantics, and task-oriented communications," *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 1, pp. 5–41, 2022.
- [3] Z. Chen, Z. Zhang, and Z. Yang, "Big AI models for 6g wireless networks: Opportunities, challenges, and research directions," *arXiv preprint arXiv:2308.06250*, 2023.
- [4] C. Chaccour, W. Saad, M. Debbah, Z. Han, and H. V. Poor, "Less data, more knowledge: Building next generation semantic communication networks," *arXiv preprint arXiv:2211.14343*, 2022.
- [5] Y. Wang, M. Chen, T. Luo, W. Saad, D. Niyato, H. V. Poor, and S. Cui, "Performance optimization for semantic communications: An attention-based reinforcement learning approach," *IEEE Journal on Selected Areas in Communications*, vol. 40, no. 9, pp. 2598–2613, 2022.
- [6] T. Han, Q. Yang, Z. Shi, S. He, and Z. Zhang, "Semantic-preserved communication system for highly efficient speech transmission," *IEEE Journal on Selected Areas in Communications*, vol. 41, no. 1, pp. 245–259, 2022.
- [7] W. Xu, Z. Yang, D. W. K. Ng, M. Levorato, Y. C. Eldar, and M. Debbah, "Edge learning for B5G networks with distributed signal processing: Semantic communication, edge computing, and wireless sensing," *IEEE J. Sel. Topics Signal Process.*, vol. 17, no. 1, pp. 9–39, Jan. 2023.
- [8] X. Peng, Z. Qin, D. Huang, X. Tao, J. Lu, G. Liu, and C. Pan, "A robust deep learning enabled semantic communication system for text," in *GLOBECOM 2022-2022 IEEE Global Communications Conference*. IEEE, 2022, pp. 2704–2709.
- [9] H. Xie, Z. Qin, G. Y. Li, and B.-H. Juang, "Deep learning enabled semantic communication systems," *IEEE Transactions on Signal Processing*, vol. 69, pp. 2663–2675, 2021.
- [10] Z. Weng and Z. Qin, "Semantic communication systems for speech transmission," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 8, pp. 2434–2444, 2021.
- [11] H. Tong, Z. Yang, S. Wang, Y. Hu, W. Saad, and C. Yin, "Federated learning based audio semantic communication over wireless networks," in *2021 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2021, pp. 1–6.
- [12] H. Xie and Z. Qin, "A lite distributed semantic communication system for internet of things," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 1, pp. 142–153, 2020.
- [13] X. Zou, "A survey on application of knowledge graph," in *Journal of Physics: Conference Series*, vol. 1487, no. 1. IOP Publishing, 2020, p. 012016.
- [14] Z. Yang, M. Chen, Z. Zhang, and C. Huang, "Energy efficient semantic communication over wireless networks with rate splitting," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 5, pp. 1484–1495, 2023.
- [15] N. Zhang, X. Xu, L. Tao, H. Yu, H. Ye, S. Qiao, X. Xie, X. Chen, Z. Li, L. Li *et al.*, "Deepke: A deep learning based knowledge extraction toolkit for knowledge base population," *arXiv preprint arXiv:2201.03335*, 2022.
- [16] C. Wang, X. Liu, Z. Chen, H. Hong, J. Tang, and D. Song, "Deepstruct: Pretraining of language models for structure prediction," *arXiv preprint arXiv:2205.10475*, 2022.
- [17] D. Fensel, U. Şimşek, K. Angele, E. Huaman, E. Kärle, O. Panasiuk, I. Toma, J. Umbrich, A. Wahler, D. Fensel *et al.*, "Introduction: what is a knowledge graph?" *Knowledge graphs: Methodology, tools and selected use cases*, pp. 1–10, 2020.
- [18] L. Floridi and M. Chiriatti, "Gpt-3: Its nature, scope, limits, and consequences," *Minds and Machines*, vol. 30, pp. 681–694, 2020.