

Breaking the Boundaries: A Unified Framework for Chinese Named Entity Recognition Across Text and Speech

Jinzhong Ning, Yuanyuan Sun*, Bo Xu, Zhihao Yang, Ling Luo, Hongfei Lin
School of Computer Science and Technology, Dalian University of Technology, China
jinzhongning1996@gmail.com
(syuan, xubo, yangzh, lingluo, hflin@dlut.edu.cn)

Abstract

In recent years, with the vast and rapidly increasing amounts of spoken and textual data, Named Entity Recognition (NER) tasks have evolved into three distinct categories, i.e., text-based NER (TNER), Speech NER (SNER) and Multimodal NER (MNER). However, existing approaches typically require designing separate models for each task, overlooking the potential connections between tasks and limiting the versatility of NER methods. To mitigate these limitations, we introduce a new task named Integrated Multimodal NER (IMNER) to break the boundaries between different modal NER tasks, enabling a unified implementation of them. To achieve this, we first design a unified data format for inputs from different modalities. Then, leveraging the pre-trained MMSpeech model as the backbone, we propose an **Integrated Multimodal Generation Framework (IMAGE)**, formulating the Chinese IMNER task as an entity-aware text generation task. Experimental results demonstrate the feasibility of our proposed IMAGE framework in the IMNER task. Our work in integrated multimodal learning in advancing the performance of NER may set up a new direction for future research in the field. Our source code is available at <https://github.com/NingJinzhong/IMAGE4IMNER>.

1 Introduction

Named Entity Recognition (NER) (Li et al., 2020a) is a fundamental and significant task in the field of Natural Language Processing and has been extensively studied to address the challenges posed by real-world text data. Chinese NER (CNER) (Liu et al., 2022), a significant subdomain of NER, specifically deals with challenges unique to Chinese, such as no clear word boundaries and the ambiguity from homophones and polyphones, drawing significant academic focus (Zhang and Yang, 2018; Li et al., 2020b; Ma et al., 2020).

*Corresponding author

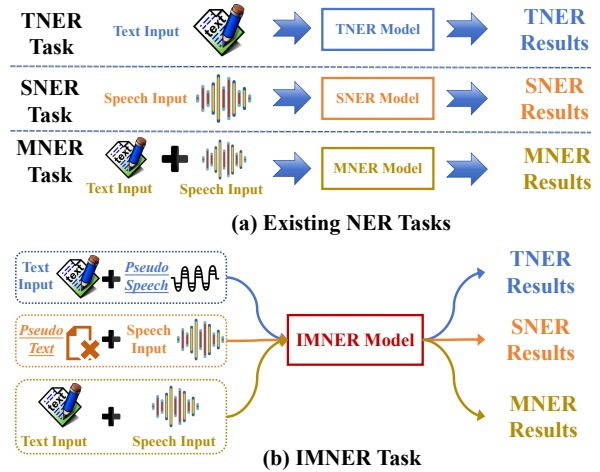


Figure 1: Comparison of existing NER tasks with the Integrated Multimodal NER (IMNER) task proposed in this paper. In the figure, TNER represents text-based NER, SNER stands for Speech NER, MNER denotes multimodal NER, Pseudo Speech refers to meaningless zero audio waveforms, and Pseudo Text indicates non-sensical text sequences.

Traditionally, Named Entity Recognition (NER) tasks have concentrated on text-based NER (TNER) (Gui et al., 2019b; Li et al., 2020b, 2022). However, as the volume of audio data increases, there has been a growing interest in Speech NER (SNER) (Yadav et al., 2020; Chen et al., 2022; Shon et al., 2022), which focuses on extracting named entities from speech, and Multimodal NER (MNER) (Sui et al., 2021; Liu et al., 2023), which involves extracting entities from both speech and text.

Currently, data on the internet often appears in multiple modalities, such as user-generated content in social media and news reports in the media, which may be in text or audio formats, or a combination of both speech and its corresponding text. The key information extracted by NER tasks, such as persons and locations, can assist in search and recommendation in the news domain, as well as in analyzing trending topics and public opinion in

social media. However, existing NER systems are usually designed for a single mode, either solely as SNER, TNER or MNER, as shown in Figure 1. These approaches face two significant issues. **Issue 1:** Treating SNER, TNER and MNER as three separate tasks overlooks the potential interconnections between them. **Issue 2:** The need to design distinct models for each of the SNER, TNER and MNER tasks limits the versatility and overall efficiency of NER methods.

Beyond the above issues, significant advances in speech processing have been achieved with Multimodal Pre-trained Models (MPMs) using data from both text and speech (Zhou et al., 2022; Ao et al., 2022). These models, trained on various tasks like speech-to-text and text-to-text generation, highlight the potential interconnections between different modal tasks. However, these MPMs face a challenge, identified here as **Issue 3:** While MPMs benefit from a unified pre-training architecture across modalities, the need for task-specific fine-tuning for various downstream applications, to some extent, restricts their universality.

To address these challenges, in this paper, we introduce a new NER task named **Integrated Multimodal NER (IMNER)**. The IMNER task aims to break the boundaries of traditional SNER, TNER and MNER by presenting a unified Named Entity Recognition (NER) framework capable of handling inputs from various modalities (text, speech, or both) to efficiently recognize Chinese named entities, as illustrated in Figure 1. Moreover, previous studies (Sui et al., 2021; Liu et al., 2023) have shown that features like pauses in speech signals can reduce ambiguities in Chinese NER tasks, which often arise from the lack of clear word delimitation or the presence of homophones. The IMNER approach, leveraging data from the SNER, TNER, and MNER tasks, possesses the potential to overcome the difficulties associated with the absence of natural word segmentation and the frequent occurrence of homophones in Chinese text.

To solve the IMNER task, our approach begins with an original design of a data format unification method that transforms the data formats of TNER, SNER and MNER tasks into a unified data scheme. As illustrated in Figure 1, we treat TNER and SNER tasks as MNER tasks with missing speech and text modalities, respectively. For these “missing” modalities, we substitute Pseudo Speech and Pseudo Text. Based on the unified data format, and using the multimodal pretrained model

MMSpeech (Zhou et al., 2022) as backbone, we introduce the **Integrated Multimodal Generation Framework (IMAGE)**, an encoder-decoder structure to execute the Chinese IMNER task. Specifically, inspired by the recent success of generative methods in NER tasks, we formulate the IMNER task as an entity-aware text generation task (Chen et al., 2022; Wang et al., 2023). Unlike previous works, our approach uniquely leverages the interrelations among the three different modalities of TNER, SNER and MNER tasks, facilitating the realization of the IMNER task.

The main contributions of this work can be summarized as follows:

- We introduce a new task, Integrated Multimodal NER (IMNER), aimed at breaking the boundaries between TNER, SNER and MNER tasks, enabling the model to uniformly handle inputs from various modalities.
- From a novel perspective, we design a unified data format for TNER, SNER and MNER, establishing a bridge between these three tasks and serving as the basis for IMNER.
- Utilizing the MMSpeech model as the backbone, we propose an Integrated Multimodal Generation Framework (IMAGE), formulating the Chinese IMNER task as an entity-aware text generation task. Notably, our IMAGE framework is capable of handling both flat and nested entity scenarios.
- Experimental results reveal that the IMAGE framework effectively exploits potential connections among TNER, SNER and MNER tasks, boosting their performance. IMAGE achieves competitive performance in these tasks, proving the viability of the IMNER task and incidentally the advantages of the IMAGE framework.

2 Related Work

2.1 Text-based Chinese NER (TNER)

In Chinese NER, the lack of natural word boundaries and the existence of homophones introduce ambiguity in the text, posing challenges for Chinese NER. Therefore, in recent years, incorporating external lexicon resources to enhance Chinese NER performance has been proven to be an effective solution and has achieved significant success (Zhang and Yang, 2018; Gui et al., 2019a,b; Li

et al., 2020b; Liu et al., 2021). Additionally, for the extraction of nested entities, recent work utilizing a unified NER framework (Li et al., 2022; Yan et al., 2021) to extract both flat and nested entities has shown promising results.

2.2 Speech NER (SNER)

Speech NER (SNER), which is essential for Spoken Language Understanding (SLU) (Caubrière et al., 2020; Shon et al., 2022), initially adopts a two-stage pipeline approach (Cohn et al., 2019): converting speech to text with Automatic Speech Recognition (ASR) and then tagging named entities in the generated text. To overcome the error accumulation inherent in this approach, End-to-End (E2E) methods for languages like French (Ghanay et al., 2018), English (Yadav et al., 2020), and Chinese (Chen et al., 2022) have emerged, which incorporate entity-aware ASR, directly integrating entity tagging into the ASR decoding process.

2.3 Multimodal NER (MNER)

With the rapid growth of multimodal data on the internet, leveraging multimodal information to enhance the performance of NER systems has attracted increasing academic attention. In the field of English NER, existing work has primarily focused on using data from text and image modalities to improve the performance of NER systems in social media contexts (Sun et al., 2021; Chen et al., 2021; Xu et al., 2022; Jia et al., 2023). Similarly, in the field of Chinese NER, Multimodal NER (MNER) that combines text with audio signals (Sui et al., 2021; Liu et al., 2023) has been introduced and achieved significant success.

2.4 MPMs Based on Text and Speech

Recently, Multimodal Pretrained Models (MPMs) have received widespread attention in the field of speech processing. In English, models such as SpeechT5 (Ao et al., 2022) and STPT (Tang et al., 2022), which propose encoder-decoder pre-training using unlabeled text and speech data, have achieved significant success. Following this, in the Chinese Automatic Speech Recognition (ASR), MM-Speech (Zhou et al., 2022) makes great improvements through a multi-modal multi-task encoder-decoder pre-training framework.

It is important to note that, in our work, although MMSpeech serves as the backbone, our model, named IMAGE, distinguishes itself from the aforementioned MPMs by overcoming the need

for individual fine-tuning across different downstream tasks. Furthermore, while recent works like SpeechGPT (Zhang et al., 2023) and Qwen-Audio (Chu et al., 2023) have demonstrated the capability to handle both speech and text inputs in conversational tasks, to our knowledge, our work is the first attempt to explore integrated modality input capability within an NER system.

3 Methodology

In this section, we first introduce the IMNER task definition. Then we detail the implementation of IMAGE, including its overall structure as depicted in Figure 2.

3.1 Task Definition

Given the input X , which may be text $\{\mathbf{x}_{\text{text}}\}$, speech $\{\mathbf{x}_{\text{speech}}\}$, or a combination of both $\{\mathbf{x}_{\text{text}}, \mathbf{x}_{\text{speech}}\}$, the goal of IMNER is to find each entity in X and then assign a label $y \in Y$, where Y is a predefined label types (e.g., PER, LOC, etc.).

3.2 Formulating IMNER into Text Generation

Inspired by the success of generative methods in TNER (Wang et al., 2023) and SNER (Chen et al., 2022), we formulating the IMNER task as an entity-aware text generation task. Illustrated by the sample in Figure 2, for the text “末阳市文联主席张三” and its associated speech waveform, the Entity-aware Text Generation Target is designated as “<(末阳市)文联>主席[张三]”. Special tokens are incorporated into the vocabulary to annotate entities in the generated text, specifically, “[]” for PER, “()” for LOC, and “< >” for ORG. We chose the entity-aware text generation task for generating entities primarily because this method allows for the simultaneous acquisition of entity span information and entity text content.

3.3 Details of the IMAGE Framework

3.3.1 Backbone Model

In this paper, we employ a multimodal pre-trained model with an encoder-decoder structure, MM-Speech (Zhou et al., 2022), as our backbone model. The original MMSpeech structure primarily consists of: (1) a multi-layer Transformer-based MM-Speech encoder shared by text and speech modalities, equipped with a multi-layer convolutional and Transformer-based speech feature extractor, and a static word vector embedding for text feature extraction; (2) a decoder composed of multiple Transformer layers.

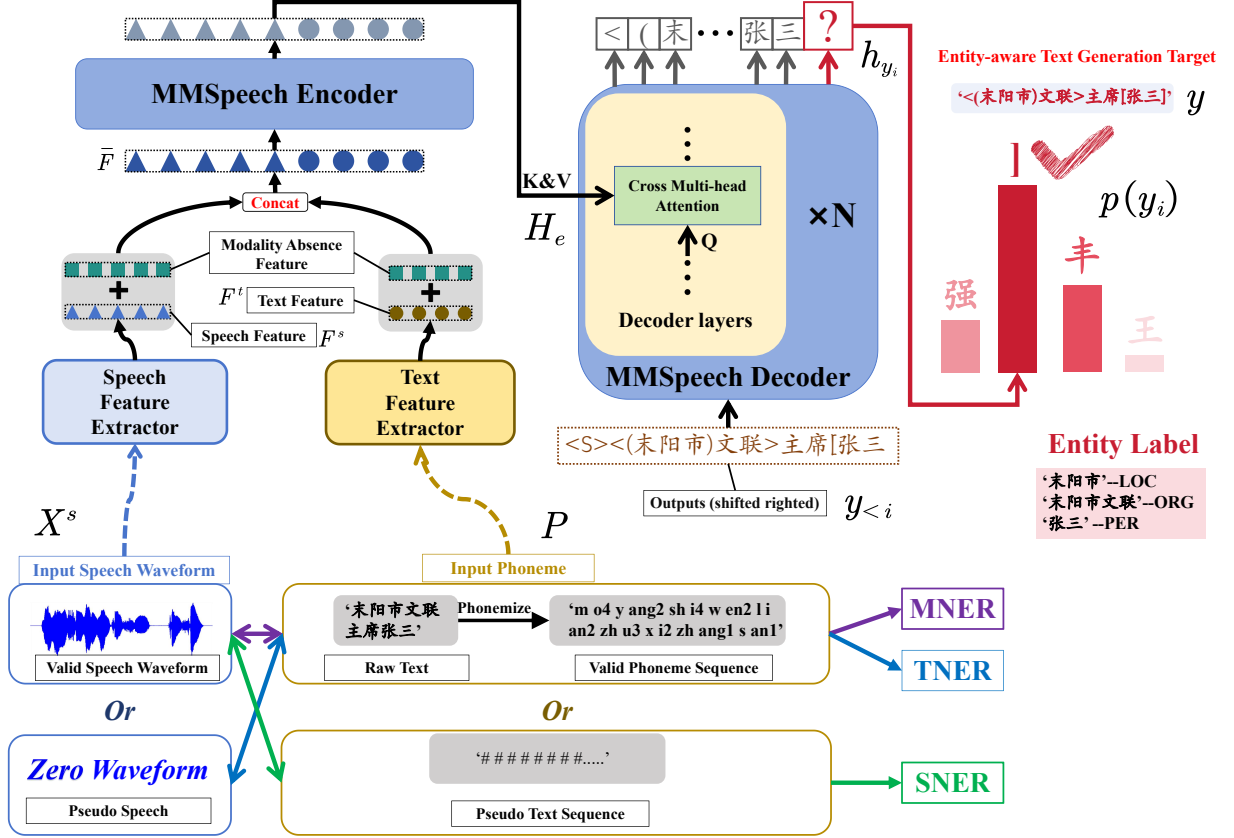


Figure 2: Overall structure of IMAGE. The figure illustrates with the example of the text “末阳市文联主席张三”(Chairman Zhang San of the Moyang City Cultural Association) and its corresponding speech waveform. In the figure, “<S>” denotes a special token indicating the start of the generated output, while “#” and “\$” respectively represent meaningless phoneme tokens and Chinese character tokens. The purple, green, and blue arrows at the bottom right of the figure explain the composition of input data for the MNER, SNER, and TNER tasks within the IMAGE framework, respectively.

Currently, MMSpeech is mainly used for downstream tasks with speech input, such as speech recognition. Besides speech data, MMSpeech was trained with a large volume (292GB) of text data, giving it a strong ability to model text. However, the capability of MMSpeech to handle text input tasks or dual input tasks with speech and text has been overlooked and not fully explored. This indicates the significant potential for expanding MMSpeech’s application across various modal scenarios.

3.3.2 Unified Integrated Modal Data Format

The IMNER task comprises three sub-tasks: SNER, TNER and MNER, each involving different modal components in the input data. To transform the integrated modal inputs of IMNER task into a unified data format, maintaining data consistency and laying the groundwork for handling all three tasks with a uniform model structure, we adopt a novel perspective within the IMAGE framework. Here,

we treat TNER and SNER tasks as MNER tasks with “missing” speech and text modalities, respectively. For these “missing” modalities, we substitute Pseudo Speech and Pseudo Text as illustrated in Figure 2.

In the Unified Integrated Modal Data Format, the Input Speech Waveform is denoted as $X^s = \{x_1^s, \dots, x_{N^s}^s\}$, where N^s represents the length of the speech waveform. When the speech modality is missing from the input, X^s represents a fixed-length sequence of all-zero signals. Because Chinese characters and their corresponding sounds are not tightly mapped to one another, the encoder in MMSpeech converts the original text input $X^t = \{x_1^t, \dots, x_{N^t}^t\}$ into phoneme input $P = \{p_1, \dots, p_{N^p}\}$, where N^t and N^p denote the sequence lengths of X^t and P , respectively.

3.3.3 Feature Extractors

Speech Feature Extractor: Consistent with MMSpeech, we first convert the raw speech waveform

into Mel-filterbank features, and then use a multi-layer convolutional network followed by a Transformer encoder (comprising multiple transformer layers with multihead self-attention (Vaswani et al., 2017)) as the speech feature extractor. The speech features F^s for the input speech X^s are computed as follows:

$$F^s = SFE(X^s) \quad (1)$$

where $SFE(\cdot)$ denotes the Speech Feature Extractor, $F^s \in \mathbb{R}^{L_{F^s} \times d_h}$, L_{F^s} is the length of the speech features, and d_h is the dimension of the hidden features (consistent with all d_h in this paper).

Text Feature Extractor: Consistent with the pre-training phase of MMSpeech, we utilize static embeddings in our model to obtain the feature representation of the input phoneme sequence P :

$$F^t = E^{(p)}(P) \quad (2)$$

where $E^{(p)}(\cdot)$ denotes the operation of static phoneme embedding, and $F^p \in \mathbb{R}^{N^p \times d_h}$.

Modality Absence Feature: To enhance the model’s ability to detect the absence of a modality, thereby encouraging it to focus on inputs from present modalities and ignore inputs from missing ones, we introduce a learned embedding to every $f_i^t (1 \leq i \leq N^p)$ and $f_i^s (1 \leq i \leq N^s)$ to incident whether that modality is missing:

$$\bar{f}_i^t = \begin{cases} f_i^t + m_{\text{missing}}, & \text{if text is missing} \\ f_i^t + m_{\text{present}}, & \text{otherwise} \end{cases} \quad (3)$$

$$\bar{f}_i^s = \begin{cases} f_i^s + m_{\text{missing}}, & \text{if speech is missing} \\ f_i^s + m_{\text{present}}, & \text{otherwise} \end{cases} \quad (4)$$

where $m_{\text{missing}} \in \mathbb{R}^{d_h}$ and $m_{\text{present}} \in \mathbb{R}^{d_h}$ are the learned embeddings indicating the absence or presence of the modality, respectively. This approach allows the model to dynamically adapt its processing based on the availability of each modality. Ultimately, we obtain the final speech feature representation $\bar{F}^s = \{\bar{f}_1^s, \dots, \bar{f}_{N^s}^s\}$ and the final text feature representation $\bar{F}^t = \{\bar{f}_1^t, \dots, \bar{f}_{N^p}^t\}$.

3.3.4 Encoder and Decoder of MMSpeech

For \bar{F}^s and \bar{F}^t , we combine them through a concatenation operation to form the feature representation \bar{F} that is fed into the MMSpeech encoder-decoder structure:

$$\bar{F} = \bar{F}^t \oplus \bar{F}^s \quad (5)$$

where $\bar{F} \in \mathbb{R}^{(N^p+N^s) \times d_h}$, and \oplus denotes the concatenation operation.

Encoder: IMAGE feeds the concatenated text and speech features representation \bar{F} into the MM-Speech encoder, which is a multi-layer Transformer encoder, to obtain the hidden representation of the integrated modal input as follows:

$$H_e = \text{Encoder}(\bar{F}) \quad (6)$$

Decoder: Afterwards, H_e is fed into the MM-Speech decoder, a multi-layer Transformer decoder, to model the probability distribution of the output text y . At the i -th step of decoding, the probability distribution $p(y_i) \in \mathbb{R}^{|V|}$ of the i -th output token y_i in y is computed as follows:

$$h_{y_i} = \text{Decoder}(H_e, y_{<i}) \quad (7)$$

$$p(y_i) = \text{Softmax}(W_{lm}h_{y_i} + b_{lm}) \quad (8)$$

where h_{y_i} is the hidden representation at the i -th decoding step, $W_{lm} \in \mathbb{R}^{|V| \times d_h}$ and $b_{lm} \in \mathbb{R}^{|V|}$ are learnable parameters in the language model (LM) head, and $|V|$ represents the size of the vocabulary.

3.3.5 Training Strategy of IMAGE

Loss Function: During the training phase, the parameters of IMAGE are optimized by minimizing the cross-entropy loss based on teacher forcing:

$$\mathcal{L} = -\frac{1}{M} \sum_{i=1}^M \sum_{k=1}^{|V|} l_{i,k} \log p(y_{i,k}) \quad (9)$$

where $l_i \in \mathbb{R}^{|V|}$ represents the ground-truth label distribution for decoding the i -th token, and M denotes the total number of tokens in the ground-truth label sentence.

Training Data Creation: In this study, we utilize MNER datasets, containing both speech and text, to create training data for TNER and SNER by artificially removing certain modality data. During training, each sample in a batch is randomly assigned as an input for MNER, SNER or TNER. These input variations are depicted in Figure 2, illustrating the method for handling inputs with varying modality presence.

4 Experiments

4.1 Dataset & Evaluation Metrics

In this study, we train and evaluate our IMAGE framework on the IMNER task using three datasets: the flat SNER dataset AISHELL-NER (Chen et al., 2022), the nested MNER dataset CNERTA (Sui

Dataset	Data Type	Entity Type Num	Sentence				Entity Type			
			train	dev	test	total	PER.	ORG.	LOC.	total
CNERTA	Text&Audio	3	34,102	4,440	4,445	42,987	8,034	12,047	16,876	36,957
AISHELL-NER	Text&Audio	3	120,098	14,326	7,176	141,600	18,642	25,351	24,611	68,604

Table 1: Statistics of the Datasets.

et al., 2021). Both the AISHELL-NER and CNERTA datasets contain Chinese text with corresponding speech, where the Chinese text is annotated with entity information. Detailed statistics of these datasets are available in Table 1. Regarding evaluation metrics, we use the F1 score (F1), commonly employed in NER tasks, to assess the model’s effectiveness.

4.2 Experimental Setting

During the training phase, we generate each training sample according to the Training Data Creation method described in Section 3.3.5. In the evaluation phase, we first manually remove the corresponding modality information from the test set to produce three versions of the test set for TNER, SNER and MNER, allowing for a comprehensive evaluation of the model’s performance.

Additionally, for the MSRA dataset, which follows the same annotation guidelines as the AISHELL-NER dataset and includes the same entity types, we employ the entire MSRA training set (46,539 samples) along with all the SNER training data from the AISHELL-NER dataset (120,098 samples) for training our model. Subsequently, we evaluate the model’s performance on the TNER task using MSRA’s test set and on the SNER task using AISHELL-NER’s test set. The rationale behind this experimental setup is to verify the effectiveness of the unified cross-modal training strategy in IMAGE on the widely used TNER dataset, i.e., MSRA.

For our IMAGE model, we initialize the model parameters using the pre-trained MMSpeech model, including components such as the Speech Feature Extractor, Text Feature Extractor, MMSpeech Encoder, and MMSpeech Decoder. We trained IMAGE using the Large¹ versions of the MMSpeech pre-trained model weights, employing teacher forcing and reported results for both. During the training phase, the training data for MNER, SNER, and TNER were balanced with a ratio of 1:1:1. We used a batch size of 12 and a learning rate

of 1e-5. During the decoding phase, we applied the beam search method with a beam width of 5. Additionally, we incorporated Pseudo Phoneme and character sequences of length 40 in our model. For the pseudo speech input, we utilized a sequence of zero values with a length of 10,000. We implemented the training of the IMAGE model using PyTorch on a GeForce RTX 4090 GPU, employing the AdamW optimizer with a warm-up rate of 0.1, and trained it for 50 epochs on each dataset.

4.3 Comparison Models

We compare the performance of several strong baseline models on the TNER, SNER and MNER tasks using the benchmark datasets employed in this paper. The baseline models used fall into three main categories: (1) methods that only use the text modality (Text-only Methods), (2) methods that only use the speech modality (Speech-only Methods), and (3) multimodal methods that use both text and speech (Multimodal Methods). The introductions to the three categories of baseline models selected for our comparison are as follows:

(1) **Text-only Methods:** We chose two types of baselines in this category. The first type includes State-of-the-Art (SOTA) methods based on Bert-large² (Cui et al., 2020), such as *Bert-large-CRF*, *FLAT* (Li et al., 2020b), and *W²NER* (Li et al., 2022). Bert-large-CRF and FLAT employ the same Nested Structure Linearization method for annotating nested entities as in the M3T (Sui et al., 2021) work. The second type consists of models with an encoder-decoder structure for NER using an entity-aware text generation task (ETG), including *Bart-large*³ (407M) (Shao et al., 2021), *MT5-base*⁴ (582M) (Xue et al., 2021), and the original version of *MMSpeech*.

(2) **Speech-only Methods:** We also chose two types of baselines in this category. The first type includes end-to-end methods based on an encoder-decoder structure, where the encoder in-

¹https://modelscope.cn/models/iic/ofa_mmspeech_pretrain_large_zh

²<https://huggingface.co/hfl/chinese-macbert-large>

³<https://huggingface.co/fnlp/bart-large-chinese>

⁴<https://huggingface.co/google/mt5-base>

Modality	Methods	CNERTA(nested)			AISHELL-NER		
		TNER	SNER	MNER	TNER	SNER	MNER
Text-only Methods	1.Bert-large-CRF(325M)	76.09 ^h	-	-	93.29 ^h	-	-
	2.FLAT(Bert-large)	79.31 ^{‡h}	-	-	93.57 ^{‡h}	-	-
	3.W ² NER(Bert-large)	79.25 [‡]	-	-	93.72[‡]	-	-
	4.Bart-large-ETG(407M)	76.84	-	-	92.82	-	-
	5.MT5-base-ETG(582M)	76.91	-	-	92.94	-	-
	6.MMSpeech-ETG(613M)	77.34	-	-	92.81	-	-
Speech-only Methods	7.Conformer-ETG(E2E)(Chen et al., 2022)	-	60.36 [‡]	-	-	73.37 [Ⓟ]	-
	8.MMSpeech-ETG(E2E)	-	69.82	-	-	75.42	-
	9.Conformer-ASR + Bert-large(Pipeline)	-	60.92	-	-	74.10	-
	10.MMSpeech-ASR + Bert-large(Pipeline)	-	69.76	-	-	74.84	-
Multimodal Methods	11.Bert-USAF(Bert-base)(Liu et al., 2023)	-	-	76.73 [Ⓟ]	-	-	-
	12.Bert-M3T(Bert-large)	-	-	79.51 ^{‡h}	-	-	93.75 ^{‡h}
	13.IMAGE(only MNER data)	-	-	79.46	-	-	93.34
IMNER Methods	14.IMAGE(MMSpeech-large,613M)	79.85	70.36	80.49	93.05	75.76	93.77

Table 2: F1-score (%) of the proposed IMAGE method and baselines on the TNER, SNER, MNER versions of the test sets for two benchmark datasets. Here, “ETG” refers to models perform NER task using an entity-aware text generation task. “(E2E)” and “(Pipeline)” respectively denote the end-to-end SNER methods and pipeline SNER methods. Superscript [‡] indicates results obtained through official implementation. Superscript [Ⓟ] denotes experimental results reported from the original paper. Superscript ^h signifies the use of the same Nested Structure Linearization method for annotating nested entities as in the M3T(Sui et al., 2021) work.

puts are speech, and the decoder annotates entities using an entity-aware text generation task, including *Conformer-ETG* (Chen et al., 2022) and *MMSpeech-ETG*. The second type involves Pipeline methods that first recognize speech into text using ASR and then annotate entities using Bert-large, such as *Conformer-ASR + Bert-large* (Chen et al., 2022), *MMSpeech-ASR*⁵ + *Bert-large*.

(3) **Multimodal Methods:** We selected the SOTA methods in Multimodal Named Entity Recognition (MNER) based on speech and text, including *Bert-USAF* (Sui et al., 2021) and *Bert-M3T* (Liu et al., 2023). Additionally, we included the *IMAGE framework trained solely on MNER data* as a baseline method.

4.4 Results and Analysis

4.4.1 Main Results

We compare our proposed IMAGE model with several strong Text-only Baselines, Speech-only Baselines, and Multimodal Baselines, with the experimental results reported in Table 2. It is evident that, unlike existing baseline models that are limited to solving single-modality tasks, IMAGE not only breaks the boundaries between modalities by simultaneously addressing TNER, SNER and MNER tasks but also achieves highly competitive performance across these tasks. From the experimental results, we can further observe that:

⁵https://modelscope.cn/models/iic/ofa_mmspeech_asr_aishell1_large_zh

	TNER	SNER	MNER
IMAGE	79.85	70.36	80.48
w/o MAF	79.27	69.84	80.03
w/o PT	79.61	69.97	80.24
w/o PS	79.43	70.08	80.17
w/o PT&PS	79.38	69.89	79.94
w/o TNER	38.75	70.12	80.34
w/o SNER	79.61	34.76	80.27
w/o MNER	79.49	70.04	41.61

Table 3: An ablation study of the IMAGE (MMSpeech-large). F1 scores (%) were evaluated on the test sets of three different tasks in CNERTA. “MAF” represents Modality Absence Feature. “PT” and “PS” respectively denote Pseudo Text Input and Pseudo Speech Input. The feature vectors corresponding to “PT” and “PS” are masked out in the Transformer through the attention mask to negate their influence. The acronyms “TNER”, “SNER”, and “MNER” specifically refer to the training data for the respective tasks.

(1) Compared to baselines using MMSpeech trained on single task data (methods 6, 8, 13 in Table 2), our proposed IMAGE method achieves significant performance improvements on TNER, SNER and MNER tasks. This demonstrates that our IMAGE method can effectively leverage the potential correlations among the three tasks across different input modalities, facilitating complementary benefits and jointly enhancing performance across all tasks.

(2) On the SNER and MNER tasks, the per-

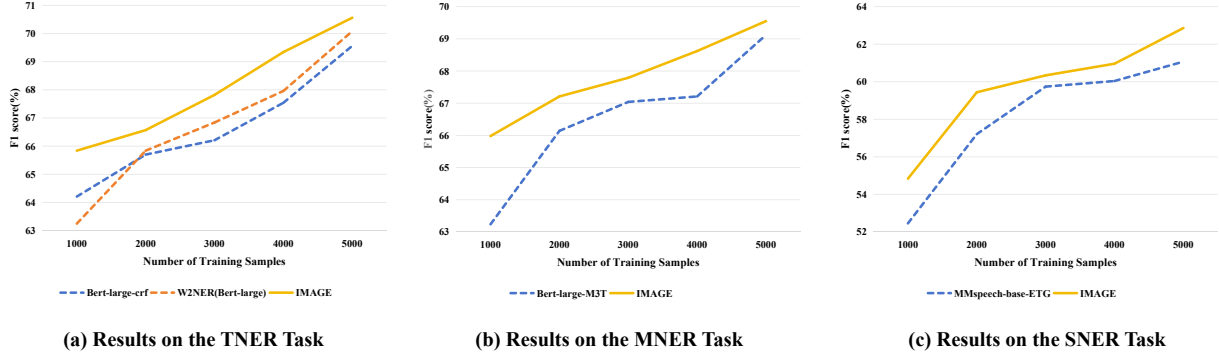


Figure 3: F1 Scores (%) of IMAGE on the three subtasks of the IMNER task in the CNERTA dataset with varying numbers of training samples. For each method depicted above, under different training data sizes, we chose identical hyperparameters and ran the experiments five times with different random seeds, averaging the F1 scores to obtain the final results.

formance of IMAGE with the MMSpeech backbone surpasses all baseline methods. This indicates that our proposed IMAGE framework can exploit the complementarity between tasks across different modalities, enhancing modeling capabilities for both speech-only and speech-text multimodal tasks.

(3) On the TNER task, our IMAGE method with MMSpeech-large backbone achieves performance comparable to the SOTA method W²NER. Moreover, the performance of the IMAGE framework exceeds all baselines based on an encoder-decoder structure using an entity-aware text generation task for entity annotation. This suggests that the IMAGE framework can improve text information encoding capabilities through joint training and modeling of tasks across different input modalities.

(4) The performance of the IMAGE method on the MNER task surpasses its performance on the TNER task, indicating that multimodal inputs combining text and speech provide more effective information than text-only data, thus enhancing model performance on NER tasks. Additionally, the SNER task not only requires entity annotation but also the accurate transcription of speech to text, increasing the complexity of the task. Therefore, the performance of the IMAGE framework on the SNER task is notably lower than on the TNER task.

4.4.2 Ablation Study

To assess the impact of various components in IMAGE, we conducted ablation experiments on the CNERTA dataset, with the findings presented in Table 3. Our conclusions are as follows:

(1) The removal of the Modality Absence Feature leads to reduced model performance, highlight-

ing its role in enhancing IMAGE’s ability to discern valid modal information in the input. Additionally, eliminating either Pseudo Text Input or Pseudo Speech Input diminishes the model’s performance. These components are believed to capture global information across different modalities, fostering synergy among the tasks within IMAGE.

(2) Excluding the training data for any one of the TNER, SNER, and MNER tasks during the training process results in a decline in the model’s performance across all three NER subtasks. This suggests that there is an intrinsic interconnection among the TNER, SNER, and MNER tasks, allowing them to mutually benefit from each other. This finding supports our IMNER approach, addressing **Issue 1** raised in the introduction of this paper, which concerns the overlooked potential interconnections between these tasks.

4.4.3 Performance on Low-resource Scenarios

We conducted experiments with training data sets of 1000, 2000, 3000, 4000, and 5000 samples to evaluate the performance of IMAGE in low-resource scenarios, with results shown in Figure 3. The experimental results reveal that, with limited training data resources, IMAGE maintains an advantage compared to baselines trained on single-task data. This demonstrates that within the low-resource context, the IMAGE framework can still effectively leverage the potential connections and complementarity among the three IMNER subtasks (i.e., TNER, SNER, MNER) to enhance the performance across these tasks. This underscores the potential of the IMAGE method in scenarios with limited training resources.

5 Conclusions

In our study, we introduce the Integrated Multimodal NER (IMNER) task, bridging the gap between text-based NER, speech NER, and multimodal NER to enable a unified approach to these three distinct tasks. By designing a novel unified data format and leveraging the pre-trained MM-Speech as backbone, we introduced the IMAGE framework, transforming the Chinese IMNER task into an entity-aware text generation task. Experimental results reveal the effectiveness of IMAGE, marking a significant step forward in integrated multimodal learning for NER, which may shed light on future research in this research domain.

Limitations

In this section, we discuss two limitations of the IMAGE framework as follows:

(1) Language Limitation: Currently, the IMAGE framework is designed to address the Chinese IMNER task exclusively. This restriction arises because the MMSpeech backbone, on which IMAGE relies, exhibits robust and balanced representation capabilities in both text and speech modalities only in Chinese. In contrast, English lacks multimodal pre-training models that perform equally well across both modalities. The available models, such as SpeechT5 (Ao et al., 2022) and STPT (Tang et al., 2022), have been pre-trained on limited text corpora, resulting in weaker text representation capabilities. Therefore, there is an urgent need to develop multimodal pre-trained models using extensive text and speech data in other languages, such as English, to support IMNER tasks in those languages.

(2) Task Limitation: At present, the IMAGE framework has only been applied to the Chinese Integrated Multimodal Named Entity Recognition (IMNER) task. Future work will involve extending the IMAGE framework to other integrated multimodal information extraction tasks. This expansion aims to fully exploit the complementary nature of different modality tasks, enhancing the overall performance and applicability of the framework.

References

Junyi Ao, Rui Wang, Long Zhou, Chengyi Wang, Shuo Ren, Yu Wu, Shujie Liu, Tom Ko, Qing Li, Yu Zhang, et al. 2022. Speecht5: Unified-modal encoder-decoder pre-training for spoken language

processing. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5723–5738.

Antoine Caubrière, Sophie Rosset, Yannick Estève, Antoine Laurent, and Emmanuel Morin. 2020. Where are we in named entity recognition from speech? In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4514–4520.

Boli Chen, Guangwei Xu, Xiaobin Wang, Pengjun Xie, Meishan Zhang, and Fei Huang. 2022. Aishellner: Named entity recognition from chinese speech. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8352–8356. IEEE.

Shuguang Chen, Gustavo Aguilar, Leonardo Neves, and Tamar Solorio. 2021. Can images help recognize entities? a study of the role of images for multimodal ner. In *Proceedings of the Seventh Workshop on Noisy User-generated Text (W-NUT 2021)*, pages 87–96.

Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. 2023. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. *arXiv preprint arXiv:2311.07919*.

Ido Cohn, Itay Laish, Genady Beryozkin, Gang Li, Izhak Shafran, Idan Szpektor, Tzvika Hartman, Avinatan Hassidim, and Yossi Matias. 2019. Audio de-identification: A new entity recognition task. In *Proceedings of NAACL-HLT*, pages 197–204.

Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Shijin Wang, and Guoping Hu. 2020. Revisiting pre-trained models for Chinese natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 657–668, Online. Association for Computational Linguistics.

Sahar Ghannay, Antoine Caubrière, Yannick Estève, Nathalie Camelin, Edwin Simonnet, Antoine Laurent, and Emmanuel Morin. 2018. End-to-end named entity and semantic concept extraction from speech. In *2018 IEEE Spoken Language Technology Workshop (SLT)*, pages 692–699. IEEE.

Tao Gui, Ruotian Ma, Qi Zhang, Lujun Zhao, Yu-Gang Jiang, and Xuanjing Huang. 2019a. Cnn-based chinese ner with lexicon rethinking. In *IJCAI*.

Tao Gui, Yicheng Zou, Qi Zhang, Minlong Peng, Jinlan Fu, Zhongyu Wei, and Xuan-Jing Huang. 2019b. A lexicon-based graph neural network for chinese ner. In *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP)*, pages 1040–1050.

Meihuizi Jia, Lei Shen, Xin Shen, Lejian Liao, Meng Chen, Xiaodong He, Zhendong Chen, and Jiaqi Li. 2023. Mner-qg: An end-to-end mrc framework

- for multimodal named entity recognition with query grounding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 8032–8040.
- Gina-Anne Levow. 2006. The third international chinese language processing bakeoff: Word segmentation and named entity recognition. In *Proceedings of the Fifth SIGHAN workshop on Chinese language processing*, pages 108–117.
- Jing Li, Aixin Sun, Jianglei Han, and Chenliang Li. 2020a. A survey on deep learning for named entity recognition. *IEEE Transactions on Knowledge and Data Engineering*, 34(1):50–70.
- Jingye Li, Hao Fei, Jiang Liu, Shengqiong Wu, Meishan Zhang, Chong Teng, Donghong Ji, and Fei Li. 2022. Unified named entity recognition as word-word relation classification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10965–10973.
- Xiaonan Li, Hang Yan, Xipeng Qiu, and Xuanjing Huang. 2020b. Flat: Chinese ner using flat-lattice transformer. In *ACL*.
- Pan Liu, Yanming Guo, Fenglei Wang, and Guohui Li. 2022. Chinese named entity recognition: The state of the art. *Neurocomputing*, 473:37–53.
- Wei Liu, Xiyan Fu, Yue Zhang, and Wenming Xiao. 2021. Lexicon enhanced chinese sequence labeling using BERT adapter. In *ACL*, pages 5847–5858.
- Ye Liu, Shaobin Huang, Rongsheng Li, Naiyu Yan, and Zhijuan Du. 2023. Usaf: Multimodal chinese named entity recognition using synthesized acoustic features. *Information Processing & Management*, 60(3):103290.
- Ruotian Ma, Minlong Peng, Qi Zhang, Zhongyu Wei, and Xuan-Jing Huang. 2020. Simplify the usage of lexicon in chinese ner. In *ACL*, pages 5951–5960.
- Yunfan Shao, Zhichao Geng, Yitao Liu, Junqi Dai, Fei Yang, Li Zhe, Hujun Bao, and Xipeng Qiu. 2021. Cpt: A pre-trained unbalanced transformer for both chinese language understanding and generation. *arXiv preprint arXiv:2109.05729*.
- Suwon Shon, Ankita Pasad, Felix Wu, Pablo Brusco, Yoav Artzi, Karen Livescu, and Kyu J Han. 2022. Slue: New benchmark tasks for spoken language understanding evaluation on natural speech. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7927–7931. IEEE.
- Dianbo Sui, Zhengkun Tian, Yubo Chen, Kang Liu, and Jun Zhao. 2021. A large-scale chinese multimodal ner dataset with speech clues. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2807–2818.
- Lin Sun, Jiquan Wang, Kai Zhang, Yindu Su, and Fangsheng Weng. 2021. Rpbert: a text-image relation propagation-based bert model for multimodal ner. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 13860–13868.
- Yun Tang, Hongyu Gong, Ning Dong, Changhan Wang, Wei-Ning Hsu, Jiatao Gu, Alexei Baevski, Xian Li, Abdelrahman Mohamed, Michael Auli, et al. 2022. Unified speech-text pre-training for speech translation and recognition. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1488–1499.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, and Guoyin Wang. 2023. Gpt-ner: Named entity recognition via large language models. *arXiv preprint arXiv:2304.10428*.
- Bo Xu, Shizhou Huang, Chaofeng Sha, and Hongya Wang. 2022. Maf: a general matching and alignment framework for multimodal named entity recognition. In *Proceedings of the fifteenth ACM international conference on web search and data mining*, pages 1215–1223.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mt5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498.
- Hemant Yadav, Sreyan Ghosh, Yi Yu, and Rajiv Ratn Shah. 2020. End-to-end named entity recognition from english speech. In *Interspeech 2020, 21st Annual Conference of the International Speech Communication Association, Virtual Event, Shanghai, China, 25-29 October 2020*, pages 4268–4272. ISCA.
- Hang Yan, Tao Gui, Junqi Dai, Qipeng Guo, Zheng Zhang, and Xipeng Qiu. 2021. A unified generative framework for various ner subtasks. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5808–5822.
- Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. 2023. [Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities](#). Preprint, arXiv:2305.11000.
- Yue Zhang and Jie Yang. 2018. Chinese ner using lattice lstm. In *ACL*, pages 1554–1564.

Xiaohuan Zhou, Jiaming Wang, Zeyu Cui, Shiliang Zhang, Zhijie Yan, Jingren Zhou, and Chang Zhou. 2022. Mmspeech: Multi-modal multi-task encoder-decoder pre-training for speech recognition. *arXiv preprint arXiv:2212.00500*.