# TRANSLINGUA: A REAL-TIME AND ACCURATE SIMULTANEOUS INTERPRETATION SYSTEM

*Chenxi Huang, Wudi Bao, Zhimu Zhang, Xin Ma*

Institute of Systems Science, National University of Singapore, Singapore 119615

## 1. BUSINESS PROBLEM

As more and more people from different countries and linguistic backgrounds interact with each other, there is a growing need to facilitate communication and understanding. The globalization of business, trade, and culture creates an increasing demand for simultaneous interpretation. In recent decades, with the development of new equipment and technology, simultaneous interpretation has improved in quality and efficiency, becoming more widespread and advanced.

Simultaneous interpretation translates audio streams into text in different languages in real-time and outputs multilingual audio content. The system allows speakers to deliver their message in their native language while the audience receives real-time interpretation in their preferred language. The primary goal of simultaneous interpretation is to facilitate communication and enable effective interaction among participants from different linguistic backgrounds. Simultaneous interpretation technology can cover a variety of industry application scenarios, such as industry forums and conferences, transnational meetings, online meetings, live programs, business public lectures, etc.

The business problem that needs to be solved in Simultaneous Interpretation Systems is the need for accurate and reliable interpretation in real-time. Finding skilled and experienced interpreters can handle the demanding requirements, such as excellent linguistic skills, knowledge of the subject matter, and the ability to think on their feet and rapid response capability. However, employing human interpreters can be expensive and inflexible for users with different demands. In contrast, an AI-powered interpretation system has the following advantages:

- It can handle multiple languages simultaneously, making them ideal for events with a large number of participants from different linguistic backgrounds.

- The system can provide a cost-effective solution by reducing the need for human interpreters, allowing businesses to allocate their budgets more efficiently.

- AI systems do not experience fatigue or make mistakes due to distraction or mishearing words. This can result in better communication and fewer misunderstandings among participants.

Therefore, using AI for simultaneous interpretation can provide businesses with significant benefits, including improved accuracy and consistency, increased scalability and flexibility, cost-effectiveness, and time-saving. These benefits can lead to better communication, increased productivity, and better outcomes, making AI-powered interpretation systems an excellent investment for businesses. Overall, our proposed system aims to achieve the following performance:

- Provide seamless communication between people who speak different languages, regardless of their location or background.

- Increase efficiency and reduce costs for businesses and organizations that regularly require interpretation services.

- Provide high-quality interpretation services that are accurate, reliable, and culturally sensitive.

- Leverage the latest advances in artificial intelligence and natural language processing to enhance the capabilities and performance of interpretation systems.

## 2. SOLUTION OVERVIEW

We integrate cutting-edge technologies such as Automatic Speech Recognition (ASR), Neural Machine Translation (NMT), and Text-to-Speech (TTS) to develop a state-of-the-art system capable of facilitating seamless translation between English and Chinese for users communicating in different languages.

The system consists of a UI interface, a speech recognition module, a translation module, and a speech generation module. When the user speaks to the UI interface, the system first recognizes the speaking language into text, then performs the text translation simultaneously and generates the speech in the target language. The following flowchart defines the specific processes of using the system:

The UI interface has a start button, and text input field together with a text output field. When the user clicks on the start button to speak, the system will recognize the speech as text and display it on the text input field. Then display the translated text on the text output field according to the current
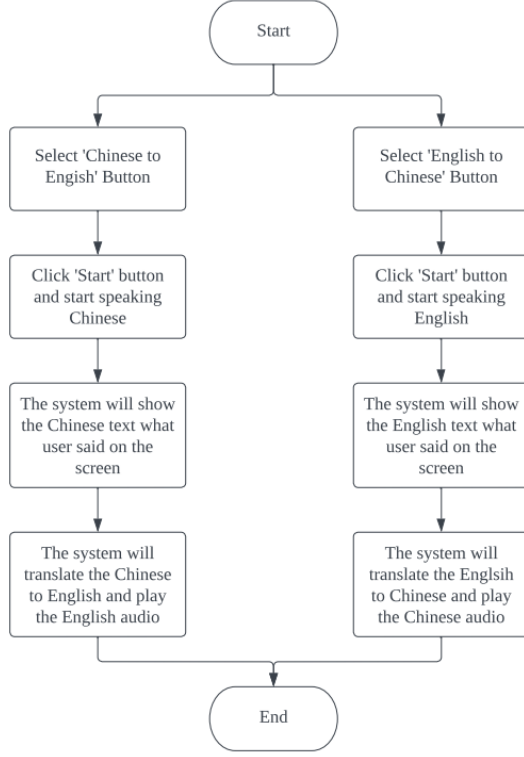
**Fig. 1**. System process

input text in real-time. Finally, it generates an audio speech out of the translated text, better with a tone of certain sentiment. Also, the user can directly input the text for translation.

We use React.js to develop the UI interface and implement Automatic Speech Recognition (ASR) in the front end. First, we set up the React application with a microphone input using the Web Speech API. Such as the getUserMedia() method that can access the user's microphone. Then we capture the user's speech input as an audio stream and use a third-party speech-to-text service to convert the audio stream into text. Finally, we display the transcribed text in our React application and send the text to the backend.

Last but not least, a good simultaneous translation system needs fluent and clear speech as the output, a well-designed text-to-speech procedure is also demanded. Although there are already sufficient TTS models out there that are well-trained to provide APIs for developers to use, we build an emotional TTS model in order to generate some more customized voices, which helps generate a more natural and contextual tone. The key contributions of our project are:

1. **Integration of Advanced Technologies for Real-Time Translation:** We've developed an innovative simultaneous translation system that seamlessly combines leading-edge technologies such as ASR, NMT, and TTS. This integration allows our system to transcribe, translate, and voice out user speech in real-time,

predicting incoming words for immediate and contextually accurate translations.

2. **Creation of an Interactive User Interface:** We've designed and built an intuitive and user-friendly UI using React.js, which encourages user interaction and fosters a seamless experience. The interface showcases live transcriptions and translations, enhancing the communication experience.

3. **Optimization of the text translation system:** We've adopted meticulous optimization methods to improve the real-time, accuracy and robustness of our text translation model. This allows it to better adapt to the scene of simultaneous translation.

## 3. PROCESS & TECHNIQUES

For the system pipeline, it first recognizes the audio speech as text, then do the real-time translation, and performs a text-to-speech task to output with voice. A certain sentiment detection method might be introduced in order to generate a speech in various tones depending on the content.

### 3.1. System Design

The whole system is developed with a frontend (React) and a backend (Django), the front end can receive the speech from the user and process it into text through an ASR module in real time. The recognized speech will be sent to the backend word by word, and then stored in a buffer, as shown in the Fig. 2.

Before proceeding to the translation module, a thread will read from the buffer to decide which couple of words can form a sentence to translate, which is also called sentence boundary detection. The reformed sentence will be popped out from the buffer and proceed to be translated.

The translated text will proceed to the TTS module, to synthesize audio from speech, then it will be encoded and sent back to the front end to display.
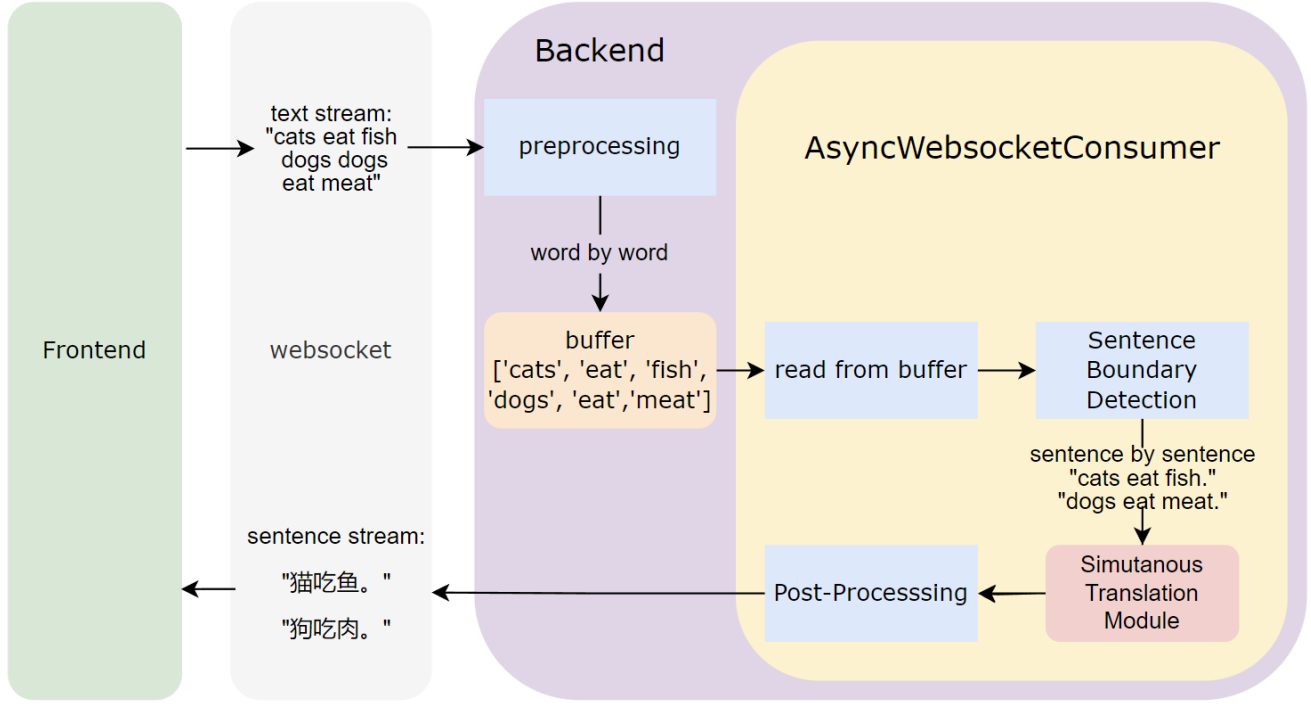
### 3.2. Automatic Speech Recognition

We use Google Cloud Speech-to-Text or Amazon Transcribe in our ASR module. These services typically offer APIs that we can use to send the audio stream and receive the transcribed text.

### 3.3. Text Translation

Our model is based on a non-autoregressive transformer (NAT) architecture.

$$p_{\text{NAT}}(Y \mid X) = \prod_{t=1}^{m} p\left(y_t \mid X\right) \tag{1}$$

**Fig. 2**. System Architecture

Unlike traditional autoregressive transformers that generate output sequentially during decoding, NAT model generates output solely based on the input, allowing for parallel decoding and significantly accelerating the inference speed, as shown in 1.

### 3.3.1. Model Architecture

The lack of interdependence between predicted words in non-autoregressive transformers may compromise translation performance. To address this issue, we adopted the GLAT[2] model, which incorporates a training strategy called *glancing*.

The GLAT model employs a two-pass decoding strategy during training:

1. In the first pass, an output prediction is obtained directly, and a portion of the target words is sampled to guide the decoder by replacing the decoder input.

2. In the second pass, the remaining unsampled words are predicted.

3. During inference, only the first pass decoding is performed to maintain fast inference speed comparable to typical non-autoregressive models.

This progressive learning strategy allows our model to capture word dependencies to the fullest extent possible. Additional optimization we adopt for the non-autoregressive transformer is the use of latent variables for decoding[1].

After obtaining the embeddings from the encoder, the model predicts the target sentence length and calculates it into the length loss. A softcopy method is adapted to predict the hidden representations of the embeddings (h) and train a latent predictor to predict the latent variables (z). The latent variables provide categorical information that helps the machine understand the tokens in different contexts. Vector quantization is applied to divide a large set of original vector representations into small groups. The hidden representations (h) and latent variables (z), both in the target dimension, are position-wise added to the input sequence (x) to predict the target language sequence (y). By combining these optimization methods, the Latent-GLAT model effectively balances the trade-off between inference speed and translation accuracy.

### 3.3.2. Optimizations

In addition to conducting research, selecting, and deploying the main architecture of the translation model, our work has made significant contributions in the following areas:

Our overall data preprocessing process is described in

1. **Data Splitting.** The bilingual Chinese-English corpus was divided by splitting using separators, and the resulting segments were stored in source and target files.
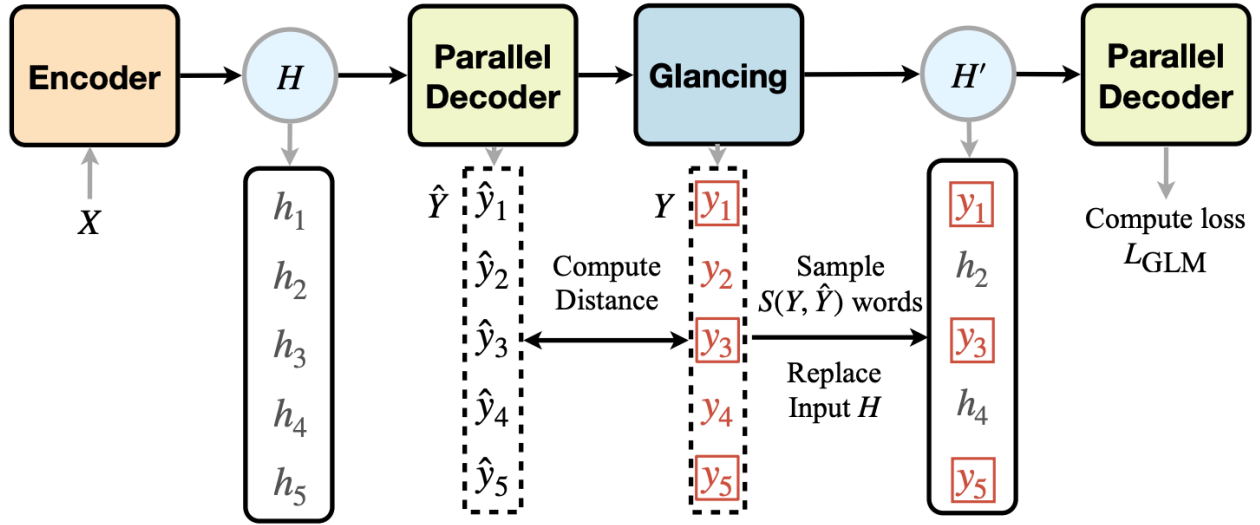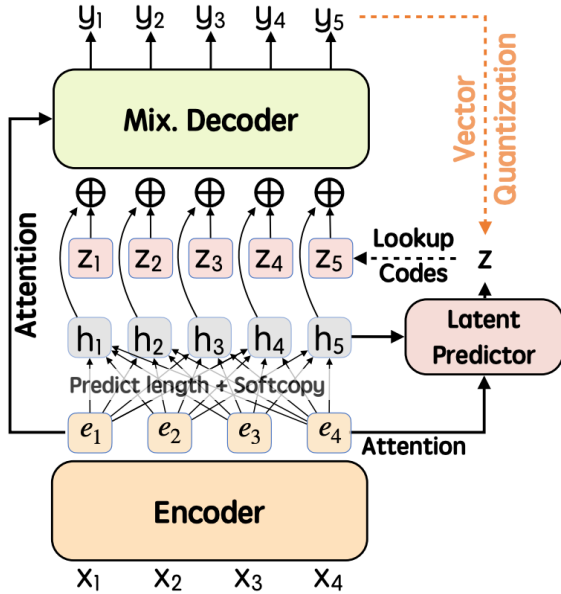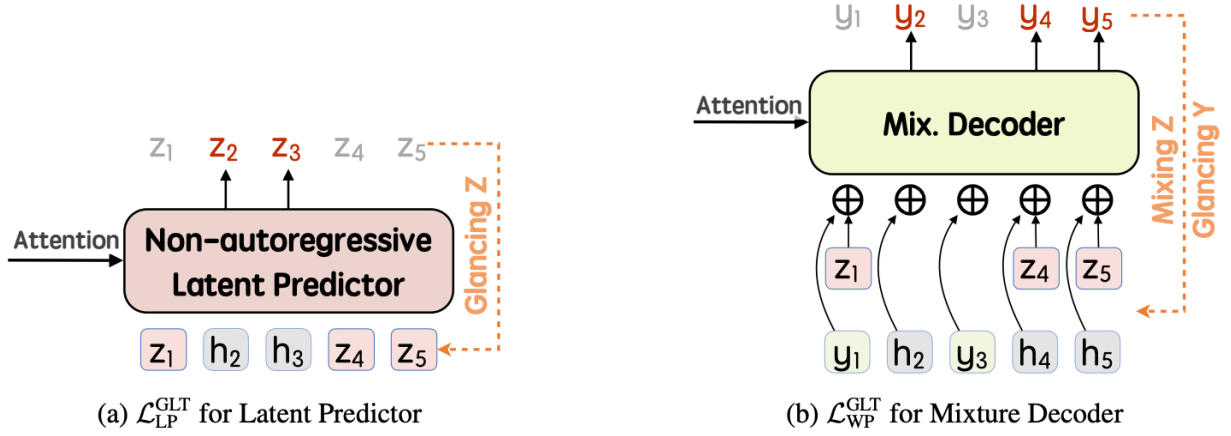
**Fig. 3**. GLAT model



**Fig. 4**. *latent*-GLAT model

2. **Normalize Punctuation.** To ensure consistent punctuation usage, we applied punctuation normalization techniques to the Chinese text in the corpus.

3. **Chinese Word Segmentation.** Employing the Jieba library, we performed Chinese word segmentation on the normalized Chinese text which had normalized punctuation. This step aimed to break down the text into individual meaningful words.

4. **Tokenization.** The preprocessed bilingual files underwent tokenization, which involved several functionalities. First, we separated English words and punctuation marks by inserting spaces between them. Next, we simplified consecutive spaces into a single space. Finally, certain symbols were replaced with their corresponding escape characters.

5. **Truecasing.** In order to address inconsistencies in capitalization, we performed case normalization on the preprocessed English text. By learning from the training data, we determined the most appropriate case forms for each English word, particularly focusing on words at the beginning of sentences.

6. **Byte Pair Encoding (BPE).** Applying BPE to the preprocessed bilingual files, we segmented the text into subword units. This approach allows the model to handle rare or out-of-vocabulary words more effectively.

7. **Data Cleaning.** To ensure high-quality data, we applied filtering techniques to the processed bilingual files, namely norm.tok.true.bpe.en and norm.seg.tok.bpe.zh. This filtering process involved removing sentence pairs

(a) $\mathcal{L}_{\mathrm{LP}}^{\mathrm{GLT}}$ for Latent Predictor  (b) $\mathcal{L}_{\mathrm{WP}}^{\mathrm{GLT}}$ for Mixture Decoder

**Fig. 5**. Training the latent predictor and mixture decoder by glancing at discrete latent variables

that fell outside a specified length range, effectively eliminating blank lines and sentence pairs with unreasonable length differences.

8. **Dataset Splitting.** Finally, the bilingual files were divided into training, testing, and development sets based on specified proportions. This ensured that the data was properly allocated for different stages of the model's training and evaluation.

Our model deployment and training process are as follows.

1. **Frequency-Sorted Dictionaries and Binary Corpus Files.** We generated dictionaries from training sets of binary corpus files. These resources, organized according to the frequency of words, facilitated efficient training and inference processes.

2. **Model Parameter Tuning.** Fairseq framework provides high flexibility to set the model. By iteratively adjusting and optimizing the parameters, we aimed to achieve the best possible translation quality on the test set.

3. **Inference Script Development.** Our project involves integrating the model's inference process with other components, necessitating the development of custom inference scripts. However, the high level of encapsulation in the fairseq framework has significantly increased our workload in this regard. Therefore, we analyze the internal workings of the fairseq framework to understand its inference pipeline and how the model generates translations. Then carefully design and implement the inference scripts. Finally, integrate it into our overall system.
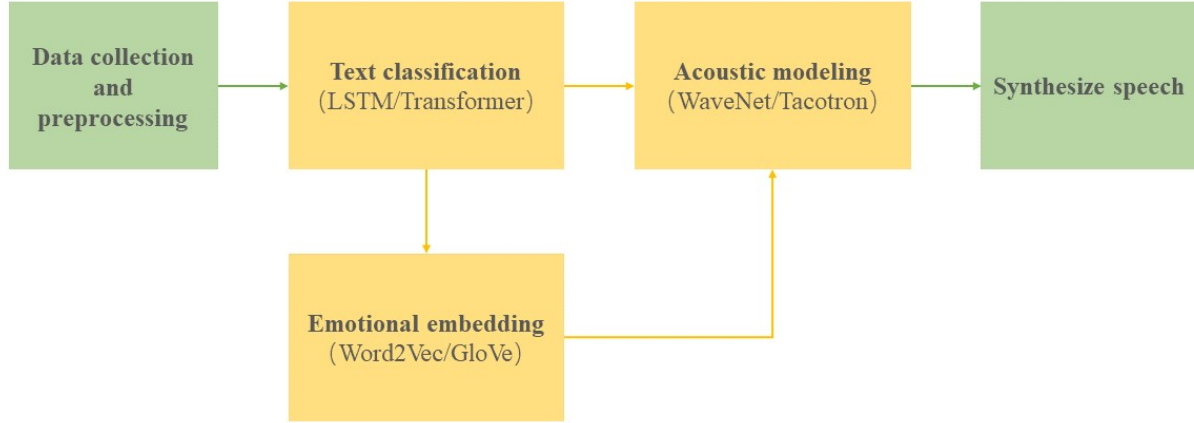
Next, we analyze the error and execute model iteration:

1. **Length Prediction Optimization.** During our experiments, we observed cases of repeated characters in the model's translations. This issue stemmed from the inherent parallelism in the model's output generation process, which led to the independence of the temporal sequence. Additionally, non-autoregressive models predict the total length of the target sentence instead of relying on a stop token for sentence completion, introducing the possibility of inaccurate length prediction and resulting in missing or overlapping characters. To address this, we introduced a progressive weight for the length prediction loss. In the later stages of training, when token prediction became relatively accurate, we gradually increased the weight of the length prediction loss. This approach encouraged the model to focus more on improving length prediction accuracy.

$$Weight = (1 - e^{-\alpha t})k \qquad (2)$$

where $Weight$ represents the weight assigned to the length prediction loss. $k$ is a scaling factor that determines the maximum weight assigned to the length prediction loss. $\alpha$ is a hyperparameter that controls the rate at which the weight increases. $t$ denotes the training progress, such as the number of training steps or epochs.

2. **Multi-branch Attention Mechanism.** We explored the utilization of a multi-branch attention mechanism[5] to enhance the model's translation capabilities. This mechanism aimed to capture different types of dependencies and relationships between words using long short-range attention, enabling the model to better understand the context and improve translation quality. We integrate this attention calculation method into our using fairseq framework. Specifically, if the number

**Fig. 6**. Emotional TTS pipeline

of branches is n, the input embedding is divided into n parts. This strategy is adopted in the multi-head self-attention layers in both encoder and decoder. The attended outputs in each branch are concatenated together following a FFN layer.

3. **Informal Dataset Fine-tuning.** Recognizing the prevalence of informal or colloquial expressions in simultaneous translation scenarios, we conducted fine-tuning on a sampled small-scale informal dataset known as the OpenSubtitles. By fine-tuning the model using this dataset, we aimed to improve its performance in handling colloquial expressions. The fine-tuning process involved selecting a diverse subset of informal expressions, merging it with the original WMT20 dataset, and using transfer learning techniques to adapt the model to informal language while retaining its proficiency in formal language.

4. **Terminology Dictionary.** To address specific domain-related terminology, we incorporated a terminology dictionary into the translation process. This dictionary contained domain-specific terms and their corresponding translations, allowing the model to accurately translate specialized vocabulary. During implementation, we annotate our translation corpus using the terminology dictionary, the dataset part. To be specific, the annotation is divided into 2 parts: TrAining Data Augmentation (TADA) and MASK[7]. With the constraints, the model can learn probabilistic replication behavior and choose whether to replace the words between $<C>$ and $</C>$ in the translation results.

### 3.4. Text-to-speech

For the TTS task, we use a general emotional TTS pipeline for our simultaneous interpretation system, as shown in Fig.

6. Finally, our emotional TTS process and techniques are described below:

1. **Data collection and preprocessing:** Collect and pre-process speech and text data, and label them with emotional tags.

2. **Text classification:** Train a text classification model (e.g., LSTM, Transformer) to classify the text into different emotional categories.

3. **Emotional embedding:** Embed the emotional tags into a low-dimensional vector space, and concatenate this emotional embedding with the linguistic embedding of the text input. This can be done using techniques like Word2Vec or GloVe.

4. **Acoustic modeling:** Train a TTS model (e.g., WaveNet, Tacotron) to generate speech from the emotional-linguistic embedding.

5. **Synthesize speech:** Use the trained model to synthesize speech with emotional expression.

To improve the accuracy of our emotion classification model, we introduce an encoder that targets raw audio instead of just using text as input. As shown in Fig. 7, first the original audio and the recognized text are encoded separately. Then, we spliced their encodings and input them into the decoder for emotional category prediction. In actual use, we use Transformer for text classification. We also use a Deep Convolutional Network with Guided Attention[3] for acoustic modeling, because according to [4], using DCN works better than using Tacotron and WaveNet on emotional TTS.

## 4. DATASETS

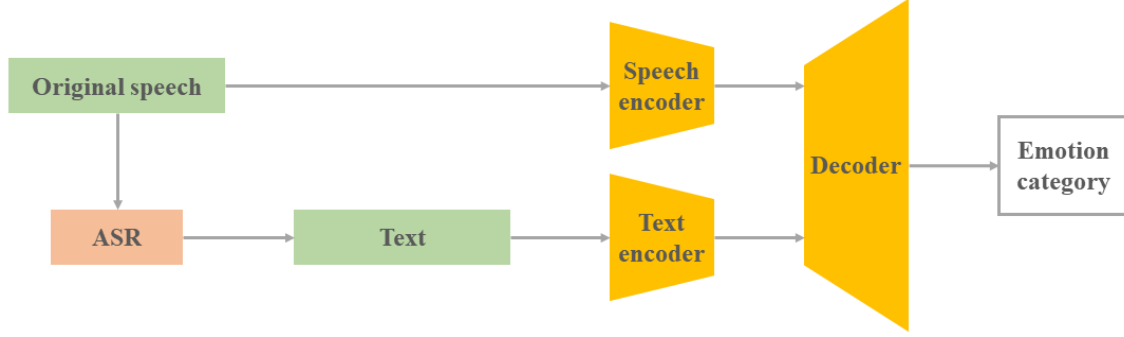In our text translation task, we utilize the WMT20 News Commentary dataset, which consists of 320K English-Chinese

**Fig. 7**. Emotional TTS optimization

| Source | People strive to o maximize their **game payoff**. |
|---|---|
| Constraint | **game payoff** -> **博弈效用** |
| TADA | People strive to o maximize their \<S\> **game payoff** \<C\> **博弈效用** \</C\>. |
| +MASK | People strive to o maximize their \<S\> **MASK MASK** \<C\> **博弈效用** \</C\>. |

**Fig. 8**. TADA and MASK

parallel sentences. This extensive collection of parallel sentences encompasses a wide range of genres, including politics, sports, entertainment, and technology. Consequently, it provides us with a vast and diverse dataset to work with.

To evaluate the performance of our translation models, we employ two essential metrics: the BLEU score for accuracy and the speedup for decoding efficiency. The BLEU score serves as a measure of similarity between the machine-generated translations and the reference translations. It quantifies how well our translations align with the expected outputs. In addition to accuracy, we also consider decoding efficiency by employing the speedup metric. Speedup refers to the decoding speed of our model compared to a traditional autoregressive transformer. This metric allows us to assess the efficiency and availability of our model for real-time simultaneous translation. By incorporating these two metrics, we ensure that our model not only delivers accurate translations but also operates efficiently in the field of simultaneous translation. In addition, to make our translation model more suitable for daily colloquial expression, we use OpenSubtitles to fine-tune our model. OpenSubtitles is a valuable resource for neural machine translation with a focus on colloquial language. It is a large-scale dataset derived from movie and TV show subtitles, including a wide range of conversations and informal expressions. With its rich collection of dialogues in multiple languages, OpenSubtitles offers an ideal opportunity

to enhance the naturalness and fluency of our translation system.

The whole OpenSubtitles dataset included over 6 million subtitles in various languages, including English and Chinese. To fine-tune our translation model, we randomly select 30K Chinese-English parallel corpus.

During the terminology enhancement, we use an open source terminology database[6]. There are approximately 2442 professional terms and 2 specialized fields, mainly focusing on basic concepts and terminology in the field of artificial intelligence.

## 5. SYSTEM PERFORMANCE

### 5.1. Experimental results

We conducted comprehensive experiments and analysis on the Chinese-English WMT20 dataset. The performance table demonstrates significant improvement in BLEU score compared to traditional non-autoregressive models, with an approximate 10-point increase. Additionally, our model achieves notable accuracy enhancements compared to GLAT models, further narrowing the performance gap with autoregressive transformers.

### 5.2. Ablation study

In order to assess the impact of individual components of the model, we conducted a series of ablation experiments to analyze their influence on translation performance and identify the most effective configurations.

We conducted an analysis of the frequency of repeated characters in the model's predictions on the test set before and after applying the progressive weight for length prediction loss. The results showed a significant decrease in the frequency of repeated characters, from 15.39% to 9.17%. This provides evidence for the effectiveness of our length prediction optimization.

The integration of the multi-branch attention mechanism improved our overall prediction BLEU score by 0.56. Addi-

**Table 1**. WMT20 results

|  | BLEU (WMT20 ZH → EN) | BLEU (WMT20 EN → ZH) | Latency(ms) | Speedups |
|---|---|---|---|---|
| Transformer (AT) | 20.80 | 18.65 | 512.3 | 1.00× |
| NAT | 7.78 | 6.36 | 33.5 | 15.29× |
| GLAT | 13.39 | 12.95 | 33.5 | 15.29× |
| Latent-GLAT | **17.14** | **15.73** | **45.3** | **11.31×** |

tionally, by reducing the dimensionality of attention calculations, the model's size was also compressed to some extent.

It is worth noting that fine-tuning the model with the informal dataset alone did not lead to an improvement in BLEU scores on our original WMT20 dataset. This is because the fine-tuning process altered the model's prediction behavior. However, through subjective evaluation, we observed a noticeable improvement in the model's adaptability to colloquial inputs. For example, it can correctly translate colloquial phrases such as "give me a hand" or "can't wait to."

By incorporating a terminology dictionary for data augmentation, the model effectively learned the replication behavior for specific domain-related terms. For inputs containing terms listed in the terminology dictionary, the model was able to generate translations that closely matched the target words specified in the dictionary.

In summary, our experiments demonstrated the effectiveness of the length prediction optimization, the performance improvement achieved through the multi-branch attention mechanism, the enhanced adaptability to colloquial inputs, and the successful utilization of the terminology dictionary for accurate term translations. These findings contribute to the overall enhancement of our machine translation system.

### 5.3. Conclusions

In conclusion, simultaneous interpretation systems play a crucial role in facilitating communication and understanding among people from different linguistic backgrounds. The globalization of business, trade, and culture has led to an increasing demand for accurate and reliable interpretation in real-time. While human interpreters have traditionally fulfilled this role, AI-powered interpretation systems offer several advantages, including the ability to handle multiple languages simultaneously, cost-effectiveness, and improved accuracy and consistency.

Our proposed solution integrates cutting-edge technologies such as Automatic Speech Recognition (ASR), Neural Machine Translation (NMT), and Text-to-Speech (TTS) to develop a novel simultaneous interpretation system. The system consists of a user-friendly UI interface, a speech recognition module, a translation module based on a non-autoregressive transformer (NAT) architecture, and a speech generation module using emotional TTS. The system allows users to speak or input text in their native language, which is then transcribed, translated, and synthesized into the target language speech in real-time.

The key contributions of our project include the integration of advanced technologies for real-time translation, the creation of an interactive user interface, the optimization of the text translation system using the GLAT model, and the development of a customizable emotional TTS model. We have also implemented various optimizations, such as length prediction optimization and multi-branch attention mechanism, to improve translation accuracy and performance. Furthermore, we have conducted data adaptation and preprocessing, model deployment and training, and error analysis and model iteration to enhance the overall system.

Overall, our simultaneous interpretation system aims to provide seamless communication between people who speak different languages, increase efficiency and reduce costs for businesses and organizations, and deliver high-quality interpretation services that are accurate, reliable, and culturally sensitive. By leveraging AI technologies, we can overcome the limitations of traditional interpretation methods and provide businesses with a valuable tool to enhance communication, productivity, and outcomes.

### 6. REFERENCES

[1] Yu Bao, Hao Zhou, Shujian Huang, Dongqi Wang, Lihua Qian, Xinyu Dai, Jiajun Chen, and Lei Li. 2022. latent-GLAT: Glancing at Latent Variables for Parallel Text Generation. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 8398–8409, Dublin, Ireland. Association for Computational Linguistics.

[2] Lihua Qian, Hao Zhou, Yu Bao, Mingxuan Wang, Lin Qiu, Weinan Zhang, Yong Yu, and Lei Li. 2021. Glancing Transformer for Non-Autoregressive Neural Machine Translation. 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 1993–2003, Online. Association for Computational Linguistics.

[3] Hideyuki Tachibana, Katsuya Uenoyama, and Shunsuke Aihara. 2018. Efficiently Trainable Text-to-Speech System Based on Deep Convolutional Networks with Guided Attention. In 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE Press, 4784–4788.

[4] Brihi Joshi, Aditya Chetan, Pulkit Madaan, Pranav Jain, Srija Anand, Eshita, Shruti Singh. (2020). An exploration into Deep Learning methods for Emotional Text-to-Speech (v1.0.0). Zenodo.

[5] Wu, Z., Liu, Z., Lin, J., Lin, Y., Han, S. (2020). Lite transformer with long-short range attention. arXiv preprint arXiv:2004.11886.

[6] https://github.com/jiqizhixin/Artificial-Intelligence-Terminology-Database

[7] Ailem, M., Liu, J., Qader, R. (2021). Encouraging neural machine translation to satisfy terminology constraints. arXiv preprint arXiv:2106.03730.