# A Comparative Study of LLM-based ASR and Whisper in Low Resource and Code Switching Scenario

*Zheshu Song, Ziyang Ma, Yifan Yang, Jianheng Zhuo, Xie Chen*[†]

MoE Key Lab of Artificial Intelligence, AI Institute
X-LANCE Lab, Department of Computer Science and Engineering,
Shanghai Jiao Tong University, Shanghai, China

{songzheshu, chenxie95}@sjtu.edu.cn

## Abstract

Large Language Models (LLMs) have showcased exceptional performance across diverse NLP tasks, and their integration with speech encoder is rapidly emerging as a dominant trend in the Automatic Speech Recognition (ASR) field. Previous works mainly concentrated on leveraging LLMs for speech recognition in English and Chinese. However, their potential for addressing speech recognition challenges in low resource settings remains underexplored. Hence, in this work, we aim to explore the capability of LLMs in low resource ASR and Mandarin-English code switching ASR. We also evaluate and compare the recognition performance of LLM-based ASR systems against Whisper model. Extensive experiments demonstrate that LLM-based ASR yields a relative gain of 12.8% over the Whisper model in low resource ASR while Whisper performs better in Mandarin-English code switching ASR. We hope that this study could shed light on ASR for low resource scenarios.

**Index Terms**: speech recognition, LLM-based ASR, Whisper, low resource, code switching

## 1. Introduction

In recent years, the scaling of model parameters has prevailed and proven effective in a range of areas, including language [1, 2], vision [3, 4], as well as speech processing [5–7]. In the realm of Automatic Speech Recognition (ASR), large-scale speech recognition models generally fall into two categories. One type is the classic end-to-end speech recognition model, exemplified by Whisper [5], which is a Transformer sequence-to-sequence model trained on various speech processing tasks and large-scale speech datasets, showing excellent performance in multilingual speech recognition and speech translation. With the advent of Large Language Models (LLMs), ASR research has increasingly shifted focus toward utilizing these models, leading to the emergence of LLM-based ASR. This approach harnesses the rich text knowledge and the reasoning ability of LLMs to improve speech recognition performance.

Currently, a wealth of impressive work [8–18] has emerged in the field of LLM-based ASR. These approaches typically employ a speech encoder network and a trainable adapter to process speech and generate embeddings, which are then passed to a decoder-only LLM. This framework seeks to strengthen the connection between acoustic features and linguistic context, enabling LLMs to better process speech input. Through a series of studies, the paradigm of enhancing speech foundation models with LLMs via projector modules has become the dominant approach in current LLM-based speech recognition research. Specifically, SALMONN [8] leverages Whisper extract semantic content and BEATs [19] for audio event information, achieving a comprehensive understanding of human speech, music, and audio events. Qwen-Audio [10] relies on Whisper as its speech encoder and enhances model performance across various audio tasks through structured task directives. SLAM-ASR [9] adopts a linear layer as the projector module achieving a new state-of-the-art performance on the 960-hour LibriSpeech [20] English task. Seed-ASR [12] employs its own powerful self-supervised audio encoder and adopts multi-stage training strategy, resulting in significant improvements over end-to-end models on comprehensive evaluation sets.

The above researches primarily focus on speech recognition for Chinese and English. Beyond these mainstream languages, low resource speech recognition remains a crucial area that cannot be overlooked. Building on the promising results of previous studies, we aim to further explore the potential of LLM-based ASR models in low resource and code switching scenario. Additionally, a comprehensive comparison between LLM-based ASR models and Whisper is conducted to evaluate their recognition performance. From our experiments, we draw the following key conclusions: (1) In languages where Whisper performs poorly, LLM-based ASR shows certain advantages. Conversely, Whisper outperforms in tasks such as Mandarin-English code switching ASR. (2) For LLM-based ASR systems, the performance of the ASR model is positively correlated with the LLM's proficiency in the specific language being recognized. We hope that our study can facilitate the research on ASR in low resource scenarios and provide valuable insights for the LLM-based ASR community.

## 2. Methods

This section primarily focuses on comparing the model architectures and underlying differences between two ASR paradigms: Whisper and LLM-based ASR.

### 2.1. Whisper

Whisper [5] is an encoder-decoder Transformer model that is capable of multiple speech tasks, including multilingual speech recognition, speech translation, language identification, and voice activity detection. The architecture of Whisper is shown in Figure 1 (a). The input to Whisper is an 80-dimensional log-Mel spectrogram of 30 seconds length $\boldsymbol{X} = [\boldsymbol{x}_1, \boldsymbol{x}_2, \cdots, \boldsymbol{x}_T]$ where T denotes the context length. The encoder blocks encode the input speech feature into hidden representations $\boldsymbol{H}$ and the decoder blocks decode the hidden representations into text tokens $\hat{\boldsymbol{y}}$ recursively conditioned on previous tokens and special prompts $\boldsymbol{p}$. In formal terms, this process can be illustrated as follows:

$$\boldsymbol{H} = AudioEncoder(\boldsymbol{X}) \tag{1}$$

$$\hat{y}_t = TextDecoder\left(p, \hat{y}_{1:t-1}, \boldsymbol{H}\right) \tag{2}$$
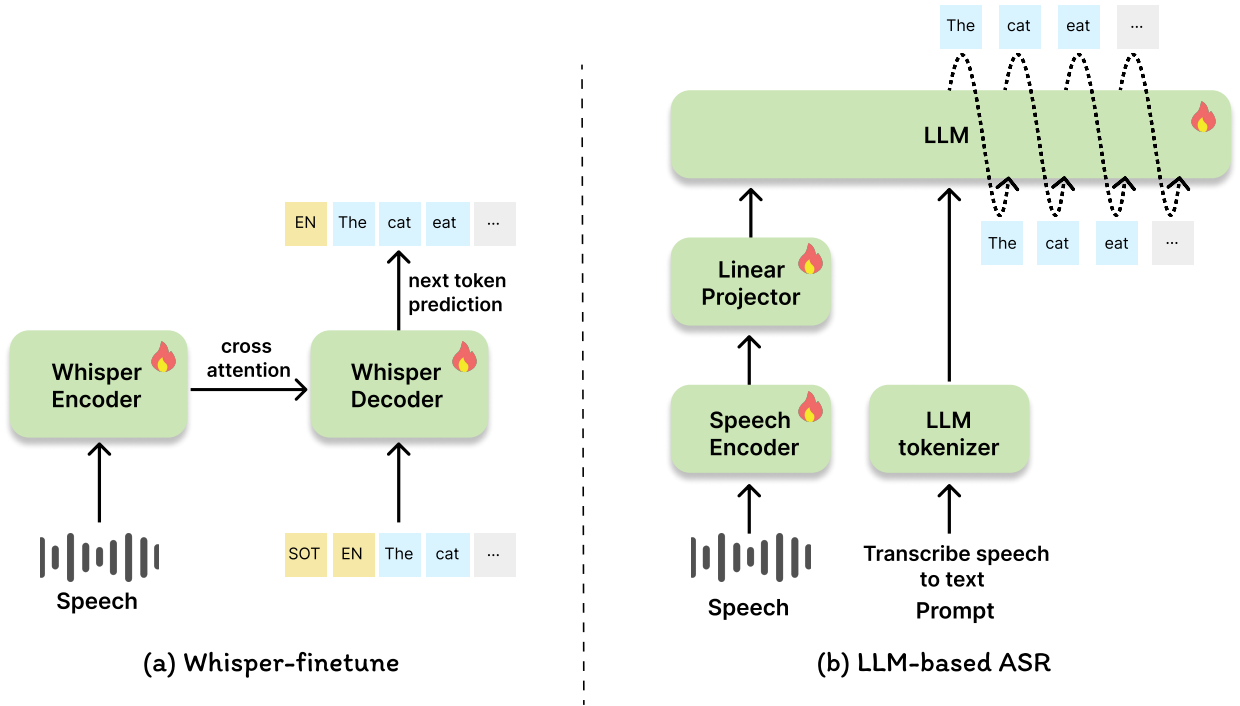
Figure 1: *Model architecture of two ASR paradigms. Left: Whisper-fientune, Right: LLM-based ASR.*

## 2.2. LLM-based ASR

As shown in Figure 1 (b), The LLM-based ASR consists of a speech encoder, a linear projector and LLM. For each sample, the given prompt (i.e., transcribe speech to text), the speech utterance, and the corresponding transcript during training are denoted as $P$, $X$, $T$, respectively.

For the input speech $X$, we first extract features by passing the speech through a speech encoder to obtain speech representations $H$, denoted as:

$$H = Encoder(X) \tag{3}$$

Due to the length of speech features is much longer than that of text features, it is necessary to downsample them. Then, the downsampled speech features are passed through a linear projector to obtain a feature sequence $E_s$ with the same dimensionality as the input to the LLM, denoted as:

$$E_s = Projector(DownSample(H)) \tag{4}$$

For the given prompt $P$ and transcription $T$, We tokenize them using the tokenizer and embedding matrix of the LLM to obtain feature sequences $E_p$ and $E_t$ as:

$$E_p = Embedding(Tokenizer(P)) \tag{5}$$

$$E_t = Embedding(Tokenizer(T)) \tag{6}$$

Next, we concatenate $E_p$, $E_s$ and $E_t$ to obtain the final feature and pass it to the LLM to obtain the output transcript $Y$, denoted as:

$$Y = LLM(E_p, E_s, E_t) \tag{7}$$

## 3. Experiments

### 3.1. Dataset

Our experiments are conducted on GigaSpeech2 dataset [21], Common Voice dataset [22] and ASRU 2019 Mandarin-English code-switching Challenge dataset [23] as shown in Table 2. In the low resource ASR experiment, we select Thai, Vietnamese, Arabic, and Welsh as representative low resource languages, drawn from GigaSpeech2 and Common Voice datasets respectively. The amount of training data for each language varies between 60 and 200 hours. For the code switching ASR experiment, approximately 200 hours of Mandarin-English data are utilized during the training phase.

Table 2: *Statistics of training and testing data (in hours)*

| Experiment | Language | Train | Test |
|---|---|---|---|
| Low-resource ASR | Thai(th) | 199.5 | 10.0 |
| | Vietnamese(vi) | 204.2 | 11.0 |
| | Arabic(ar) | 63.71 | 12.7 |
| | Welsh(cy) | 104.2 | 8.3 |
| Code-Switching ASR | Mandarin-English | 200.0 | 20.4 |

### 3.2. Training configuration

The architecture of LLM-based ASR consists of three main components: a speech encoder, a projector, and a large language model (LLM). In our experiments, the Whisper encoder serves as the speech encoder, while a linear projector is employed to align the outputs of the speech encoder with the inputs of the LLM. For the LLM component, specific models [24–27] are se-

Table 1: *Comparison of Whisper and LLM-based ASR in WER/CER for low resource languages.*

| Model | Lauguage | | | | Average WER |
|---|---|---|---|---|---|
| | **th** | **vi** | **ar** | **cy** | |
| Whisper | 20.44 | 17.94 | 15.92 | 31.86 | 21.54 |
| Whisper-finetune | 17.08 | 16.15 | **12.22** | 18.06 | 15.88 |
| LLM-based ASR | **16.13** | **15.10** | 13.30 | **10.88** | **13.85** |

lected based on the target language, as detailed in Table 3. For example, the Sailor model [24], which is fine-tuned specifically for Southeast Asian languages, is utilized for Thai and Vietnamese speech recognition.

Table 3: *Configuration of two components in LLM-based ASR.*

| Language | Speech Encoder | Large Language Model |
|---|---|---|
| th | | Sailor-7b |
| vi | Whisper Encoder | Sailor-7b |
| ar | | Jais-7b |
| cy | | Mixtral-7b |
| cn-en | | Qwen2.5-7b |

For Whisper-finetune, LoRA adaptation [28] is applied to Whisper large-v3 following LoRA-Whisper [29]. In specific, low-rank matrices where rank $r = 16$ are added to the attention layer $\{\boldsymbol{W}_k, \boldsymbol{W}_q, \boldsymbol{W}_v\}$ and fully-connected layer $\boldsymbol{W}_{fc}$ in each transformer layer in both encoder and decoder. In the training stage, we fix all the parameters of Whisper and optimize the LoRA modules with AdamW [30] with a peak learning rate of 1e-5. The number of training epochs is set to 10. All models are trained with 4 NVIDIA RTX 3090 24GB GPUs. In the testing stage, beam search with $beamsize = 5$ is employed to decode the test set.

For LLM-based ASR, we fix the parameters of the first 30 layers of the Whisper encoder, while allowing the parameters of the last two layers to remain trainable. As for LLM, low-rank matrices where rank $r = 24$ are added to the attention layer $\{\boldsymbol{W}_q, \boldsymbol{W}_v\}$. Similar to the Whisper-finetune approach, only the trainable parameters are optimized using AdamW [30] with a peak learning rate of 1e-4. The number of training epochs is also set to 10. All models are trained with 2 NVIDIA A800 80GB GPUs. Beam search with $beamsize = 5$ is applied during the testing stage.

### 3.3. Results and analysis

#### 3.3.1. Low resource ASR

As can be seen in Table 1, among the four low resource languages, the LLM-based ASR outperforms both Whisper and Whisper-finetune in three of them (Thai, Vietnamese, and Welsh), with average Word Error Rate (WER) reduced by 35.7% and 12.8% compared to Whisper and Whisper-finetune, respectively. Specifically, Whisper exhibits limited performance in Welsh recognition, with WER of 31.86% on the Welsh test set, which improves to 18.06% after LoRA fine-tuning. In contrast, LLM-based ASR demonstrates significantly better performance. This is due to Mixtral's incorporation of rich textual information and its strong support for European lan-

guages, resulting in notably higher recognition accuracy on Welsh compared to Whisper. As for Arabic, the performance of LLM-based ASR is better than Whisper, but slightly inferior to Whisper-finetune. This may be attributed to the Jais large language model's relatively limited support for Arabic, as well as the comparatively small amount of Arabic training data available.

#### 3.3.2. Code Switching ASR

From Table 4, it can be seen that LLM-based ASR performs well in the Mandarin-English code switching task, with MER reduced to below 8%, which is a 20% relative improvement compared to Whisper. Despite that, the performance of LLM-based ASR is still far behind that of Whisper-finetune. This is because Whisper demonstrates outstanding recognition performance in high resource languages such as Chinese and English, and is capable of handling Mandarin-English code-switching task with ease after fine-tuning. Therefore, it is a better choice to use Whisper's fine-tuned model to handle such tasks.

Table 4: *Comparison of Whisper and LLM-based ASR in Mandarin-English code switching ASR. MER denotes mixed error rate for both Chinese character and English words.*

| Model | CN CER | EN WER | MER |
|---|---|---|---|
| Whisper | 7.85 | 35.18 | 10.01 |
| Whisper-finetune | **4.42** | **22.39** | **6.12** |
| LLM-based ASR(Ours) | 5.97 | 26.84 | 7.99 |
| LLM-based ASR(Others) [14] | 5.13 | 29.36 | 7.76 |

### 3.4. Ablation Study

In LLM-based ASR, specific LLM is chosen based on the target language. To validate the effectiveness of this approach, a series of experiments have been conducted on Welsh and Vietnamese language.

Specifically, the recognition performance of Welsh and Vietnamese is evaluated using four large language models: Sailor-7b, Mixtral-7b, Llama3-8b, and Vicuna-7b. The Sailor model has been specifically fine-tuned for Southeast Asian languages, making it particularly well-suited for tasks involving these languages. In contrast, Mixtral is pretrained on a substantial corpus of European languages, thus offering enhanced performance for European language tasks. Experimental results are shown in Table 5. It is evident that the model achieves optimal performance when the large language model is closely aligned with the target language, demonstrating the significance of choosing an appropriate LLM based on the target language.

Table 5: *Ablation study of LLM-based ASR on different LLMs*

| Language | Performance on different LLMs | | | |
|---|---|---|---|---|
| | Sailor-7b | Mixtral-7b | Llama3-8b | Vicuna-7b |
| cy | 11.32 | **10.88** | 11.74 | 11.19 |
| vi | **15.10** | 17.11 | 15.62 | 17.95 |

# 4. Conclusion

This study presents a comparative analysis of LLM-based ASR and Whisper in low resource and code switching scenarios. Based on extensive experiments, the following key conclusions are drawn: (1) In languages where Whisper exhibits limited performance, LLM-based ASR demonstrates certain advantages. In contrast, Whisper excels in tasks such as Mandarin-English code-switching ASR. (2) The performance of LLM-based ASR model is positively correlated with the proficiency of the LLM in the specific language being recognized. We hope that our study can facilitate the research on ASR for low resource scenarios.

# 5. References

[1] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess *et al.*, "Scaling laws for neural language models," *arXiv preprint arXiv:2001.08361*, 2020.

[2] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai *et al.*, "Training compute-optimal large language models," *arXiv preprint arXiv:2203.15556*, 2022.

[3] J. Betker, G. Goh, L. Jing, TimBrooks, J. Wang *et al.*, "Improving image generation with better captions," *Computer Science*, vol. 2, 2023.

[4] M. Dehghani, J. Djolonga, B. Mustafa, P. Padlewski, J. Heek *et al.*, "Scaling vision transformers to 22 billion parameters," in *Proc. ICML*, Honolulu, 2023.

[5] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey *et al.*, "Robust speech recognition via large-scale weak supervision," in *Proc. ICML*, Honolulu, 2023.

[6] Y. Zhang, W. Han, J. Qin, Y. Wang, A. Bapna *et al.*, "Google USM: scaling automatic speech recognition beyond 100 languages," *arXiv preprint arXiv:2303.01037*, 2023.

[7] V. Pratap, A. Tjandra, B. Shi, P. Tomasello, A. Babu *et al.*, "Scaling speech technology to 1,000+ languages," *Journal of Machine Learning Research*, vol. 25, 2024.

[8] C. Tang, W. Yu, G. Sun, X. Chen, T. Tan *et al.*, "SALMONN: Towards generic hearing abilities for large language models," in *Proc. ICLR*, Vienna, 2024.

[9] Z. Ma, G. Yang, Y. Yang, Z. Gao, J. Wang *et al.*, "An embarrassingly simple approach for LLM with strong ASR capacity," in *arXiv preprint arXiv:2402.08846*, 2024.

[10] Y. Chu, J. Xu, X. Zhou, Q. Yang, S. Zhang *et al.*, "Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models," in *arXiv preprint arXiv:2311.07919*, 2023.

[11] Y. Chu, J. Xu, Q. Yang, H. Wei, X. Wei *et al.*, "Qwen2-audio technical report," in *arXiv preprint arXiv:2407.10759*, 2024.

[12] Y. Bai, J. Chen, J. Chen, W. Chen, Z. Chen *et al.*, "Seed-ASR: Understanding diverse speech and contexts with LLM-based speech recognition," in *arXiv preprint arXiv:2407.04675*, 2024.

[13] H. Xue, W. Ren, X. Geng, K. Wei, L. Li *et al.*, "Ideal-LLM: Integrating dual encoders and language-adapted LLM for multilingual speech-to-text," in *arXiv preprint arXiv:2409.11214*, 2024.

[14] F. Zhang, W. Geng, H. Huang, C. Yi, and H. Qu, "Boosting code-switching ASR with mixture of experts enhanced speech-conditioned LLM," in *arXiv preprint arXiv:2409.15905*, 2024.

[15] W. Yu, C. Tang, G. Sun, X. Chen, T. Tan *et al.*, "Connecting speech encoder and large language model for ASR," in *Proc. ICASSP*, Seoul, 2024.

[16] X. Geng, T. Xu, K. Wei, B. Mu, H. Xue *et al.*, "Unveiling the potential of LLM-based ASR on Chinese open-source datasets," in *arXiv preprint arXiv:2405.02132*, 2024.

[17] G. Yang, Z. Ma, F. Yu, Z. Gao, S. Zhang *et al.*, "MaLa-ASR: Multimedia-assisted LLM-based ASR," in *Proc. Interspeech*, Kos Island, 2024.

[18] Y. Gong, H. Luo, A. H. Liu, L. Karlinsky, and J. R. Glass, "Listen, think, and understand," in *Proc. ICLR*, Vienna, 2024.

[19] S. Chen, Y. Wu, C. Wang, S. Liu, D. Tompkins *et al.*, "BEATs: Audio pre-training with acoustic tokenizers," in *arXiv preprint arXiv:2212.09058*, 2022.

[20] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: An ASR corpus based on public domain audio books," in *Proc. ICASSP*, Brisbane, 2015.

[21] Y. Yang, Z. Song, J. Zhuo, M. Cui, J. Li *et al.*, "Gigaspeech 2: An evolving, large-scale and multi-domain ASR corpus for low-resource languages with automated crawling, transcription and refinement," in *arXiv preprint arXiv:2406.11546*, 2024.

[22] R. Ardila, M. Branson, K. Davis, M. Henretty, M. Kohler *et al.*, "Common voice: A massively-multilingual speech corpus," in *arXiv preprint arXiv:1912.06670*, 2020.

[23] X. Shi, Q. Feng, and L. Xie, "The ASRU 2019 Mandarin-English code-switching speech recognition challenge: Open datasets, tracks, methods and results," in *arXiv preprint arXiv:2007.05916*, 2020.

[24] L. Dou, Q. Liu, G. Zeng, J. Guo, J. Zhou *et al.*, "Sailor: Open language models for south-east asia," in *arXiv preprint arXiv:2404.03608*, 2024.

[25] A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary *et al.*, "Mixtral of experts," in *arXiv preprint arXiv:2401.04088*, 2024.

[26] N. Sengupta, S. K. Sahu, B. Jia, S. Katipomu, H. Li *et al.*, "Jais and jais-chat: Arabic-centric foundation and instruction-tuned open generative large language models," in *arXiv preprint arXiv:2308.16149*, 2023.

[27] Qwen Team, "Qwen2.5: A party of foundation models," September 2024. [Online]. Available: https://qwenlm.github.io/blog/qwen2.5/

[28] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu *et al.*, "LoRA: Low-rank adaptation of large language models," in *arXiv preprint arXiv:2106.09685*, 2021.

[29] Z. Song, J. Zhuo, Y. Yang, Z. Ma, S. Zhang *et al.*, "Lora-whisper: Parameter-efficient and extensible multilingual ASR," in *Proc. Interspeech*, Kos Island, 2024.

[30] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," in *arXiv preprint arXiv:1711.05101*, 2019.