

Soniox Chinese Speech Recognition Benchmarks

Soniox Inc, October 2023

support@soniox.com

<https://soniox.com>

Abstract

Soniox has conducted an extensive evaluation of the accuracy of different speech recognition providers in the industry. The benchmarks are summarized as follows:

- **Providers evaluated:** Soniox, OpenAI, Google, AWS, Azure, NVIDIA Riva, Deepgram, AssemblyAI, Speechmatics, and Rev AI.
- **Languages evaluated:** Chinese (simplified characters).
- **Evaluation datasets:** real-world dataset of YouTube videos varying in acoustic conditions, speaking styles, accents, topics, and the number of speakers.
- **Ground truth transcriptions** were transcribed and double-reviewed by humans, then normalized to ensure a fair evaluation across different providers.
- **Processing modes evaluated:** asynchronous transcription (file or batch) and streaming (live) transcription.
- **Results:**
 - Soniox achieved the most accurate speech recognition results in both async and streaming modes.
 - Second place was obtained by Azure; however, there was an extremely large accuracy gap between Soniox and Azure. **Soniox achieved a 61% higher accuracy** rate than Azure on YouTube videos.
 - The lowest accuracy was observed with Google, NVIDIA, and Rev AI.

Results

Character Error Rates (CER)

For each dataset, we evaluated all providers in **asynchronous** (file or batch) and **streaming** (live) modes. The metric reported is Character Error Rate (CER), which measures the percentage of characters misrecognized. A lower CER indicates a more accurate speech recognition system.

Dataset	Hours	Mode	Soniox	OpenAI	Google	AWS	Azure	NVIDIA	Deepgram	AssemblyAI	Speechmatics	Rev AI
YouTube videos	10	Async	4.64%	21.12%	31.61%	20.97%	8.79%	31.38%	22.12%	21.10%	12.82%	22.24%
		Streaming	5.62%	N/A	26.81%	21.61%	8.79%	14.42%	22.84%	N/A	12.82%	21.61%

Soniox achieved the lowest CER (highest accuracy) in both modes. The second place belongs to Azure, however, the difference in accuracy between Soniox and Azure is very large, considering this a character error rate and not a word error rate. All other providers are further behind Soniox.

Note that OpenAI (Whisper) and AssemblyAI are only supported in Async mode and not in Streaming mode.

Soniox Improvement

The table below shows how much of an improvement Soniox speech recognition AI brings over other providers, if you were to transcribe the same audio.

Dataset	Hours	Mode	Soniox	OpenAI	Google	AWS	Azure	NVIDIA	Deepgram	AssemblyAI	Speechmatics	Rev AI
YouTube videos	10	Async	N/A	86%	91%	86%	61%	91%	87%	86%	76%	87%
		Streaming	N/A	N/A	85%	82%	47%	71%	83%	N/A	67%	82%

Soniox outperforms other providers by extremely large margins. For example, in Async mode, **Soniox makes 61% fewer errors than Azure**. That is, Soniox fixes 61% of the errors that Azure's speech recognition AI makes.

Note, the improvement metric has been computed by taking the original CER reported in the table above, from which we subtracted the estimated amount of error in ground truth transcriptions. We were conservative in these estimates and estimated only 2% of human error for the YouTube videos dataset.

Methodology

Datasets

To represent **real-world** speech recognition use cases, we selected our evaluation dataset to include difficult real-world speech environments, conversations with crosstalk and interjections. All datasets are spoken in the Chinese language.

Dataset	Difficulty Level	Example	Duration
YouTube videos	Medium / High	Documentary TV Show Product Review	10 hours

Metric

We calculated the Character Error Rate (CER) following the standard definition:

$$CER = \frac{\text{Number of Characters Recognized Incorrectly}}{\text{Number of Characters in Ground Truth}}$$

$$\text{Number of Characters Recognized Incorrectly} = \text{Substitutions} + \text{Insertions} + \text{Deletions}$$

Ground Truth Transcriptions

1. Ground truth transcriptions were labeled and double-reviewed by humans.
2. Ground truth and provider transcriptions were then normalized to achieve a fair evaluation across different providers. See examples below.

Normalization Type	Normalization Performed
Punctuations	Punctuations were removed
Capitalization	Catarina ⇒ catarina

Models Evaluated

We evaluated the latest and most accurate models from each provider in both streaming and asynchronous processing modes. To do so, we integrated with every provider's API carefully following their documentation. All models were evaluated on September 28th - 29th, 2023.

Provider	Async Model	Streaming Model
Soniox	zh_v2	zh_v2 lowlatency
OpenAI Whisper	large-v2 (zh)	Streaming is not supported
Google	default (zh)	default (zh)
AWS	Amazon Transcribe (zh-CN)	Amazon Transcribe (zh-CN)
Azure	Standard (zh-CN)	Standard (zh-CN)
NVIDIA Riva	Riva ASR (zh-CN)	Riva ASR (zh-CN)
Deepgram	General Base (zh-CN)	General Base (zh-CN)
AssemblyAI	Asynchronous Transcription (zh)	Streaming is not supported
Speechmatics	Transcribe a File Enhanced (cmn)	Transcribe in Real-time Enhanced (cmn)
Rev AI	Asynchronous Speech-to-Text API (cmn)	Streaming Speech-to-Text API (cmn)

Contact Information

Email: support@soniox.com

If you have any questions, comments or suggestions about the benchmarks, feel free to reach out to us. We welcome your feedback and are always looking for ways to improve the benchmarks.

If you would like us to include a new evaluation dataset in the next benchmarking report, also reach out to us and we would be happy to do so.

Benchmark Your Dataset

If you would like to evaluate speech recognition providers on your dataset, please reach out to us, and we will assess the providers of interest on your dataset and provide you with the complete results, including the outputs of every provider.