

Soniox Spanish 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:** Spanish.
- **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.
 - Speechmatics, AssemblyAI, and Azure were among the top contenders, but they all fell far behind Soniox. For example, in streaming mode, **Soniox achieved a 55% higher accuracy** rate than Azure on YouTube videos.
 - The lowest accuracy was observed with Google, NVIDIA, and Rev AI.

Results

Word Error Rates (WER)

For each dataset, we evaluated all providers in **asynchronous** (file or batch) and **streaming** (live) modes. The metric reported is Word Error Rate (WER), an industry standard metric for evaluating the accuracy of speech recognition systems, which measures the percentage of words misrecognized. A lower WER 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 | 4.13% | 7.59% | 10.81% | 8.72% | 7.80% | 11.54% | 9.76% | 5.93% | 6.00% | 10.10% |
| | | Streaming | 4.64% | N/A | 10.20% | 15.05% | 7.83% | 12.69% | 13.83% | N/A | 7.81% | 12.59% |

Soniox achieved the lowest WER (highest accuracy) in both modes. The second place belongs to Speechmatics, AssemblyAI and Azure, however, the difference in accuracy between Soniox and any second-place provider is large. All other providers fall far behind Soniox and incur much higher WERs.

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 | 62% | 76% | 68% | 63% | 78% | 73% | 46% | 47% | 74% |
| | | Streaming | N/A | N/A | 68% | 80% | 55% | 75% | 78% | N/A | 55% | 75% |

Soniox outperforms other providers by extremely large margins. For example, in Streaming mode, **Soniox makes 55% fewer errors than Azure.** That is, Soniox fixes 55% of the errors that Azure's speech recognition AI makes.

Note, the improvement metric has been computed by taking the original WER 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 Spanish language.

| Dataset | Difficulty Level | Example | Duration |
|----------------|------------------|---|----------|
| YouTube videos | Medium / High | Panel Discussion Various Conversations | 10 hours |

Metric

We calculated the Word Error Rate (WER) following the standard definition:

$$WER = \frac{\text{Number of Words Recognized Incorrectly}}{\text{Number of Words in Ground Truth}}$$

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

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 |
|--------------------|--------------------------------------|
| Filler words | Filler words (e.g. uhm) were removed |
| 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 October 9th, 2023.

| Provider | Async Model | Streaming Model |
|-----------------|--|---|
| Soniox | es_v2 | es_v2_lowlatency |
| OpenAI Whisper | large-v2 (es) | Streaming is not supported |
| Google | latest_long (es-ES) | latest_long (es-ES) |
| AWS | Amazon Transcribe (es-ES) | Amazon Transcribe (es-US) |
| Azure | Standard (es-ES) | Standard (es-ES) |
| NVIDIA Riva | Riva ASR (es-ES) | Riva ASR (es-ES) |
| Deepgram | General Nova (es) | General Nova (es) |
| AssemblyAI | Asynchronous Transcription (es) | Streaming is not supported |
| Speechmatics | Transcribe a File Enhanced (es) | Transcribe in Real-time Enhanced (es) |
| Rev AI | Asynchronous Speech-to-Text API (es) | Streaming Speech-to-Text API (es) |

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.