Soniox Speech-to-Text Benchmarks

Soniox Inc, March 2025 support@soniox.com https://soniox.com

Abstract

Soniox conducted a comprehensive evaluation of the accuracy of various speech recognition providers in the industry. The benchmarking results are summarized as follows:

- **Providers evaluated:** Soniox, OpenAI, Google, AWS, Azure, NVIDIA, Deepgram, AssemblyAI, Speechmatics, and ElevenLabs.
- Languages evaluated: 60 languages.
- Evaluation datasets: Real-world datasets of YouTube videos for each language, covering diverse acoustic conditions, speaking styles, accents, topics, and speaker variations.
- **Ground truth transcriptions:** Transcribed and double-reviewed by humans, then normalized to ensure a fair evaluation across providers.
- Processing modes evaluated: Asynchronous (file/batch) transcription.
- **Results:** Soniox achieved the highest speech recognition accuracy across most languages by a significant margin.

Results

Evaluated Providers and Languages

To assess the accuracy of speech recognition providers across multiple languages, we conducted a rigorous benchmarking study using **Word Error Rate (WER)** and **Character Error Rate (CER)** as the primary evaluation metrics. These industry-standard metrics provide a quantitative measure of transcription accuracy, with lower values indicating superior performance. CER was used for the following languages: Korean, Chinese, Japanese, and Thai.

Our evaluation was based on **45 to 70 minutes of real-world audio per language**, sourced from YouTube to ensure a diverse and challenging dataset. The selected samples encompass **various acoustic conditions**, **speaking styles**, **accents**, **and topics**, providing a robust assessment of model performance in real-world scenarios. **Ground truth transcriptions** were carefully transcribed and double-reviewed by humans, then normalized to ensure consistency and fairness across all providers.

The table below presents accuracy results (WER and CER) for each speech recognition provider across the evaluated languages. An empty cell in the table indicates that the provider does not support that language.

Comparative Performance Analysis

In addition to tabulated results, we present visual comparisons of Soniox's performance against other providers. These plots illustrate accuracy differences between Soniox and competitors such as OpenAI, Google, AWS, and others across all languages.

Methodology

Evaluation Process

- 1. **Dataset Selection:** For each language, real-world YouTube videos were chosen to reflect diverse speech patterns, varying acoustic conditions, and multiple speaker types.
- 2. **Transcription & Ground Truth Creation:** All dataset transcriptions were manually created, double-reviewed by humans, and normalized to provide a consistent reference for evaluation. The normalization process included removing punctuation and ignoring capitalization; otherwise, the ground truth transcription remained unchanged.
- 3. **Model Integration:** Each provider's API was carefully integrated according to official documentation, ensuring a fair and accurate comparison.
- 4. Evaluation Metrics:
 - WER (Word Error Rate): Measures transcription errors at the word level.
 - CER (Character Error Rate): Used for logographic or non-space-separated languages to provide a finer-grained accuracy measurement.
- 5. **Processing Mode:** All models were evaluated in **asynchronous (file/batch) transcription mode** to ensure consistency in testing.

Models Evaluated

Provider	Model Evaluated
Soniox	stt-async-preview
OpenAl	Whisper large-v3
Google	long (for supported languages), chirp_2 (for other languages)
AWS	Best/Default
Azure	Best/Default

Provider	Model Evaluated						
NVIDIA	<pre>conformer-{lang}-asr-offline-asr-bls-ensemble (for supported languages), parakeet-1.1b-unified-ml-cs-universal-multi-asr-offline -asr-bls-ensemble (for other languages)</pre>						
Deepgram	nova-3 (for English), nova-2 (for other languages)						
AssemblyAl	best (for supported languages), nano (for other languages)						
Speechmatics	enhanced						
ElevenLabs	scribe_v1						

This evaluation provides a transparent and rigorous comparison of speech recognition performance across industry-leading providers for 60 languages.

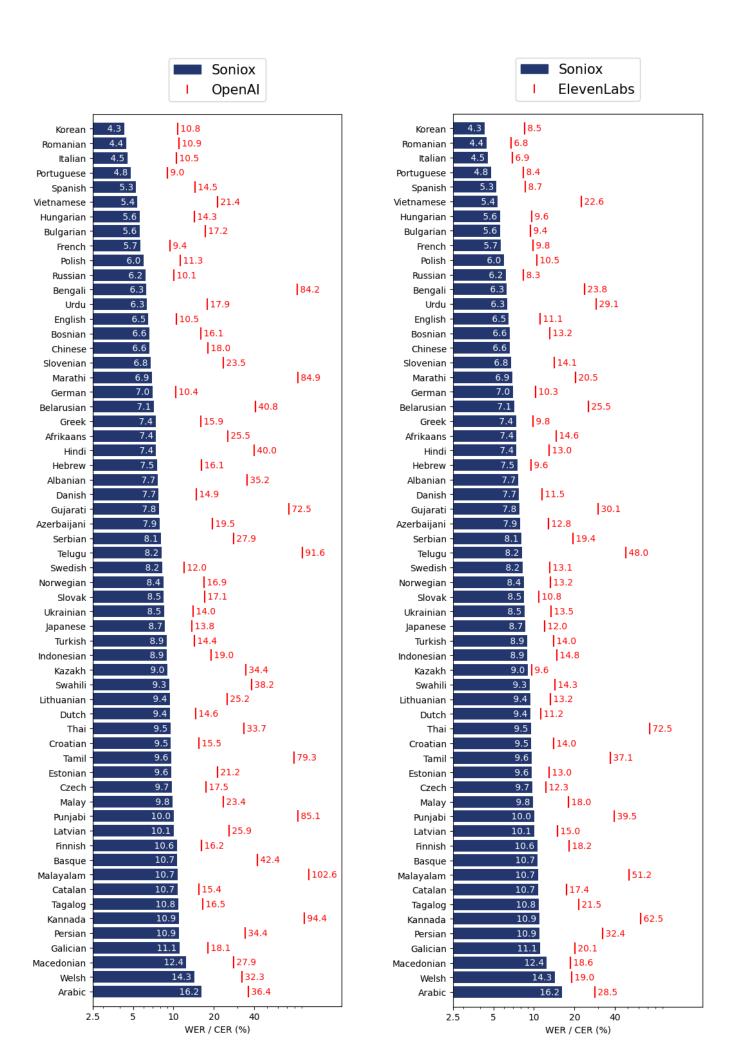
Table: Evaluated Providers and Languages

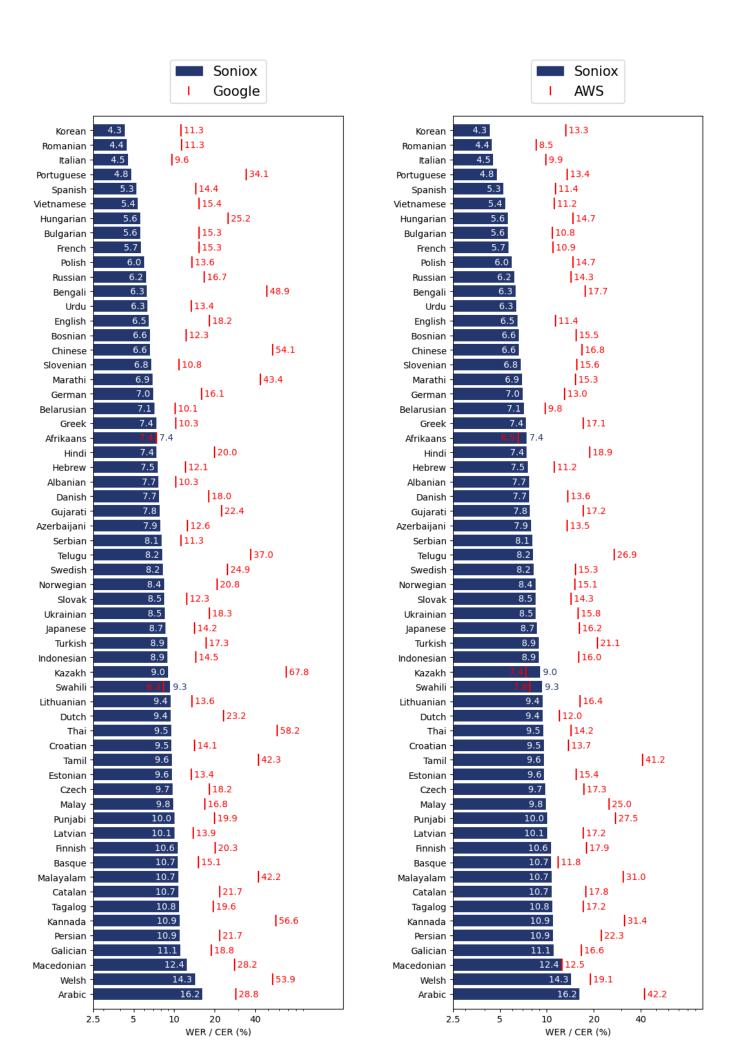
Language	Soniox	OpenAl	Google	AWS		Speech- matics	AssemblyAl	Deepgram	NVIDIA	ElevenLabs
Afrikaans	7.4	25.5	7.4	6.5			64.0			14.6
Albanian	7.7	35.2	10.3		7.9		91.1			
Arabic	16.2	36.4	28.8	42.2	37.1	22.9	55.6		44.5	28.5
Azerbaijani	7.9	19.5	12.6	13.5	19.9		49.5			12.8
Basque	10.7	42.4	15.1	11.8	18.2	7.2	77.4			
Belarusian	7.1	40.8	10.1	9.8		10.2	74.1			25.5
Bengali	6.3	84.2	48.9	17.7	40.7	27.7	101.1			23.8
Bosnian	6.6	16.1	12.3	15.5	19.6		34.2			13.2
Bulgarian	5.6	17.2	15.3	10.8	13.9	9.2	36.2	24.8	3	9.4
Catalan	10.7	15.4	21.7	17.8	24.6	13.9	34.7	19.3	3	17.4
Chinese	6.6	18.0	54.1	16.8	10.1	14.2	19.9	94.4	17.8	•
Croatian	9.5	15.5	14.1	13.7	21.1	13.4	38.2			14.0
Czech	9.7									
Danish	7.7			13.6						
Dutch	9.4									
English Estonian	6.5 9.6					9.3 12.3				11.1 13.0
Finnish	10.6									18.2
French	5.7									
Galician	11.1								10.1	20.1
German	7.0			13.0					16.5	
Greek	7.4				13.2					9.8
Gujarati	7.8						101.7			30.1
Hebrew	7.5		12.1	11.2					87.2	
Hindi	7.4									
Hungarian	5.6									9.6
Indonesian	8.9	19.0	14.5	16.0	17.0	15.5	28.7	17.1		14.8
Italian	4.5	10.5	9.6	9.9	9.0	7.4	7.6	10.8	12.2	6.9
Japanese	8.7	13.8	14.2	16.2	14.0	10.3	14.8	11.7	21.9	12.0
Kannada	10.9	94.4	56.6	31.4	58.8		100.4			62.5
Kazakh	9.0	34.4	67.8	7.4	13.9		81.5			9.6
Korean	4.3									
Latvian	10.1									15.0
Lithuanian	9.4								6	13.2
Macedonian							57.1			18.6
Malay	9.8								2	18.0
Malayalam	10.7						102.8			51.2
Marathi	6.9									20.5
Norwegian	8.4								90.5	
Persian	10.9					24.2				32.4
Polish	6.0									
Portuguese Punjabi	4.8			13.4 27.5			8.0 100.9		18.3	8.4 39.5
Romanian	4.4								5	6.8
Russian	6.2									
Serbian	8.1				*	11.7	36.6		, 10.7	19.4
Slovak	8.5				16.2	11.8			3	10.8
										14.1
Slovenian	6.8									

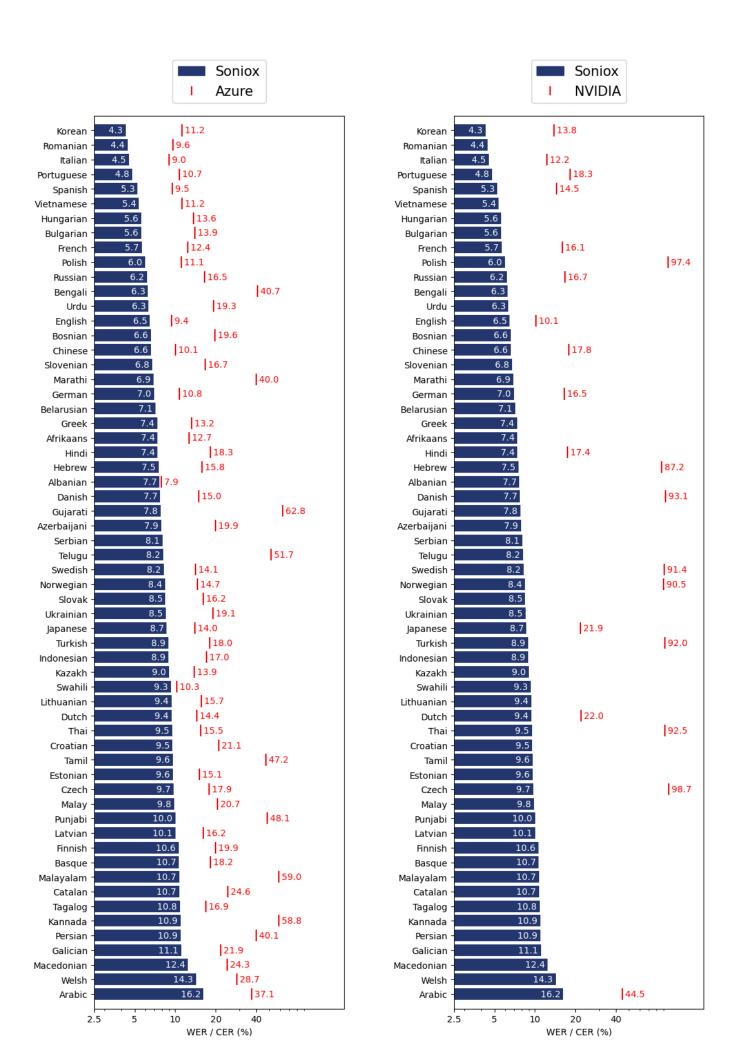
Spanish	5.3	14.5	14.4	11.4	9.5	7.1	7.1	10.1	14.5	8.7
Swahili	9.3	38.2	8.3	7.8	10.3	11.4	84.5			14.3
Swedish	8.2	12.0	24.9	15.3	14.1	11.2	28.4	13.6	91.4	13.1
Tagalog	10.8	16.5	19.6	17.2	16.9		34.5			21.5
Tamil	9.6	79.3	42.3	41.2	47.2	37.6	88.5			37.1
Telugu	8.2	91.6	37.0	26.9	51.7		99.9			48.0
Thai	9.5	33.7	58.2	14.2	15.5	13.3	85.6	48.6	92.5	72.5
Turkish	8.9	14.4	17.3	21.1	18.0	13.3	11.6	15.4	92.0	14.0
Ukrainian	8.5	14.0	18.3	15.8	19.1	15.7	13.3	14.1		13.5
Urdu	6.3	17.9	13.4		19.3	33.9	42.6			29.1
Vietnamese	5.4	21.4	15.4	11.2	11.2	8.9	12.1	14.2		22.6
Welsh	14.3	32.3	53.9	19.1	28.7	19.7	75.8			19.0

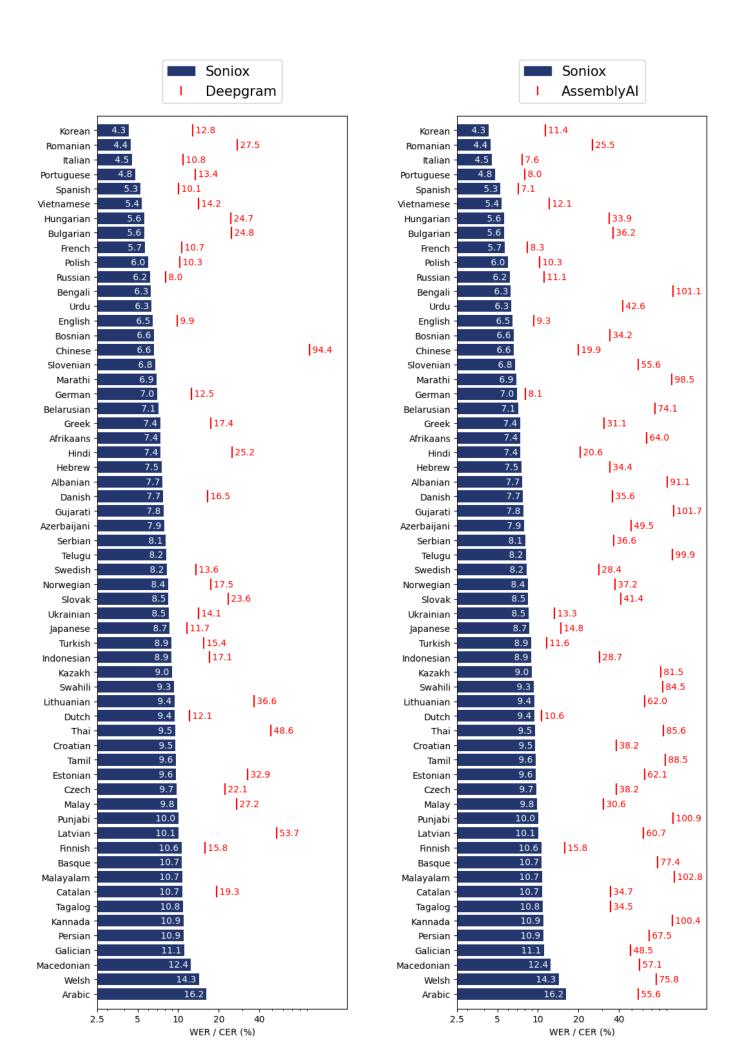
Notes:

- * For the Serbian language, AWS and Azure transcribe in Cyrillic, not Latin, whereas the ground truth transcripts are in Latin; therefore, the WER results are omitted.
- We manually reviewed the returned transcripts of provider-language pairs with high WER (e.g., AssemblyAl on Bengali) and found that the high WER was due to the provider's STT model generating low-quality transcripts rather than other technical reasons.









Soniox Speechmatics

