

End-to-End 音声基盤モデルの開発・事業成果概要

実施者

株式会社Kotoba Technologies Japan

事業概要

End-to-Endで音声を汎用的に入力・出力し処理することのできる音声基盤モデルを1.データの準備、2.基盤モデルの開発、3.学習モデル評価のステップに分けて行う。

1.データセットの構築: インターネット上から50万時間分の高品質、日本語・英語混合データを収集する。

2.基盤モデルの開発: 上記のデータセットを活用して、{1.3B, 7B}の音声基盤モデルの分散並列学習を行う。下記の評価セットで特に日本語のドメインで最先端の性能を達成する。最終的に7Bのモデル学習を成功させ、TTS Arenaなどで汎用性において日本語最高性能を目指す。

3.学習モデル評価構築: 音声基盤モデルの評価を整え、日本語音声下流評価タスクを新たに作成しコミュニティに貢献する。東北大学と連携して、今後の基盤音声モデル開発の礎となる評価データセットを作る。

社会実装イメージ

当事業で開発された音声基盤モデルは、APIやB2B音声アプリケーションとして、社会実装をすすめる。本モデルは、リアルタイムで音声出力から入力まで一貫して行えるモデルであることを重視している。これが実現されれば、コールセンター、広告、教育、福祉、医療、エンターテイメントなど、音声を用いるアプリケーションにおいて、大規模な音声AIの利活用が期待できる。

※本資料に掲載する製品名等は、各社の商標または登録商標です。

事業成果

1.データセットの構築: 日英合算で50万時間のハイクオリティ訓練データを構築した。特に日本語においては、従来にないスケールで、データ収集及びクリーニングのパイプラインを確立し、HuggingFace DatasetsなどのAI分野のスタンダードとなるツールも駆使しながら、データ構築の基盤技術を完成させた。

2.基盤モデルの開発: 1.3B, 7Bのモデルの訓練は完了し、7Bモデルにおいては国内外の最大級の日英音声基盤モデルとなり、TTS Arena、S2S Translation Arena、音声Retrieval、音声Perplexityにおいても最高ベースラインを遥かに上回る性能を記録した。日本語音声生成においてOpenAIやGoogle TTSを上回る評価を獲得した。また日英の音声翻訳においても非常に優れた精度を発揮してデモビデオとしてシステムを公開した。

3.学習モデル評価構築: 日本語音声において世界初の試みであるTTS Arena Japaneseをリリースし、1200件以上の評価ポイントが集まっている。また、コールセンターを意識したSpeaker Clusteringデータセットも構築した。

Japanese TTS Arena: Benchmarking Japanese TTS Models in the Wild

Vote to help the community find the best available text-to-speech model!

This arena is inspired and built on [TTS Arena](#).

We are actively maintaining this project. Suggestions via contact/discussion are welcome!

[Vote](#) [Leaderboard](#) [About](#)

Leaderboard

Vote to help the community determine the best Japanese text-to-speech (TTS) models.

The leaderboard displays models in descending order of how natural they sound (based on votes cast by the community).

Important: In order to help keep results fair, the leaderboard hides results by default until the number of votes passes a threshold. Tick the [Reveal preliminary results](#) to show models without sufficient votes. Please note that preliminary results may be inaccurate.

order	name	score	votes
#1	KOTOBALO-SPEECH-SPK4	1267	140
#2	KOTOBALO-SPEECH-SPK1	1246	112
#3	BLANE-TTS	1226	141
#4	OPENAI-TTS	1224	100
#5	MOE-VITS	1218	133
#6	KOTOBALO-SPEECH-SPK3	1203	110
#7	KOTOBALO-SPEECH-SPK2	1196	133
#8	GOOGLE-TTS	1154	107
#9	AMITARO-VITS	1151	115
#10	BARK	1093	121

Show all models, including models with very few human ratings.
 Reveal preliminary results

Refresh

Citation