

## INTRODUCTION

- Modern singing voice synthesis (SVS) voice-banks can sing cross-lingually.
- SVS and music generation applications have no feature for automatic song translation.
- Potential of cross-lingual SVS is untapped without automatic song translation.
- Vocal and lingual processing methods have more focus and greater performance in speech than singing.
- No *complete system* for transcription, singable translation and synthesis faithful to the original song exists.
- We present the first singing-voice to singing-voice translation (SV2SVT) system.

## PIPELINE BREAKDOWN

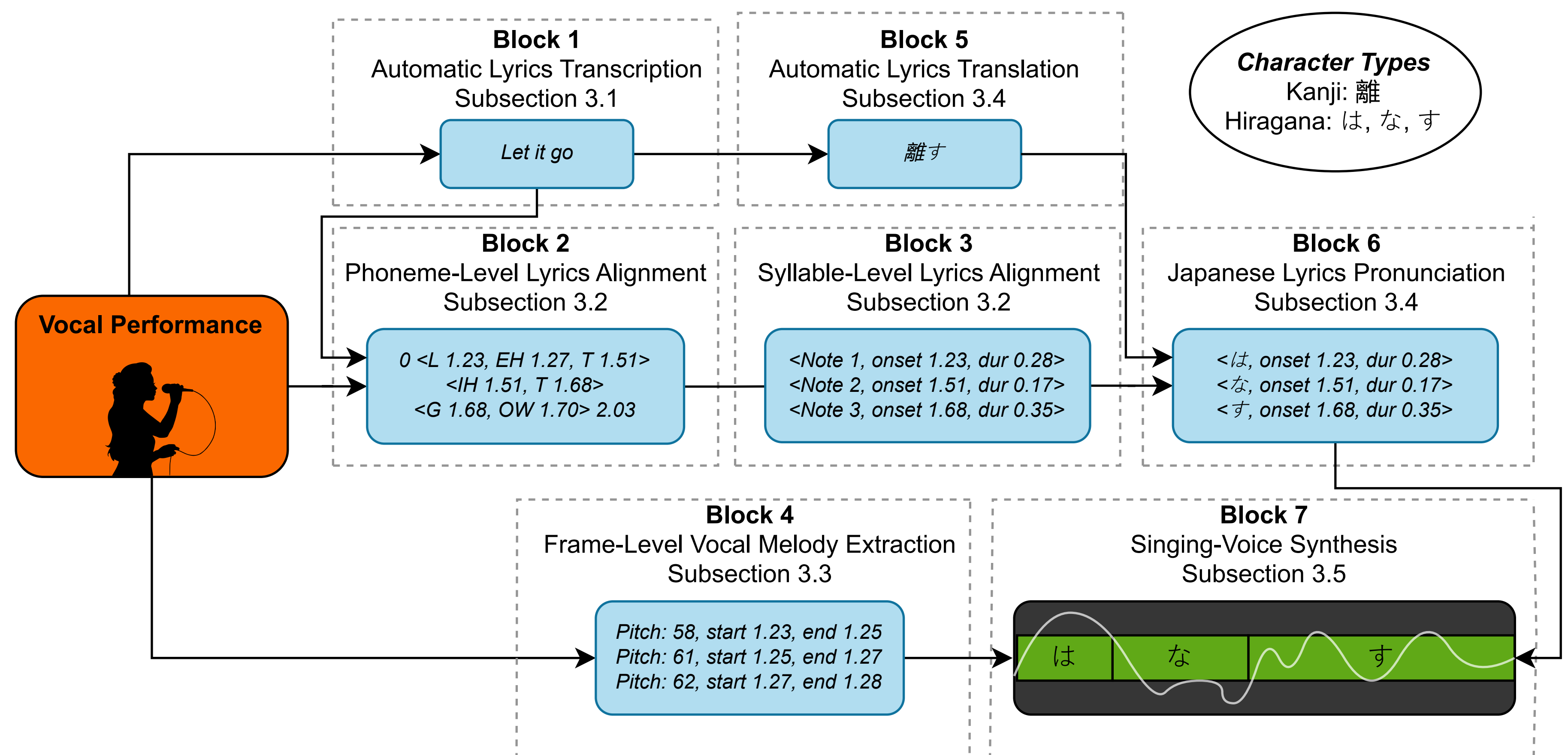
Record English vocal performance.

- **Block 1:** Transcribe lyrics.
- **Block 2:** Align phonetic sequence of the transcribed lyrics to audio.
- **Block 3:** Categorize phoneme sequences into syllables with onsets and durations.
- **Block 4:** Frame-level vocal-contour transcription.
- **Block 5:** Generate a list of possible Japanese translations.
- **Block 6:**
  - Break down all kanji in all translated sequences into their hiragana-form pronunciations and separate the words (Japanese uses no whitespace).
  - Choose the translated sequence with the most appropriate amount of syllables as best translation.
  - Align Japanese syllables (mora) with prior onsets and durations of English syllables.
- **Block 7:** Define notes in SVS engine with Japanese syllables bounded by onsets and durations. Automate the pitch values of notes with vocal-contour.

## CONCLUSION

- We have presented the first SV2SVT system in the literature.
- Our method facilitates a baseline framework for future development in English to Japanese SV2SVT as well as a baseline methodology for other language combinations.
- Further research is required for robust SV2SVT. Alternative technologies and methods must be tested, and robust tools must be developed. This will primarily be our future work.

## PROPOSED METHOD



## EXPERIMENTS

Method	Module
Whisper-Large-V3	Lyrics transcription ( <i>Block 1</i> )
Schufu lyrics-aligner	Phoneme-level lyrics alignment ( <i>Block 2</i> )
CMU Pronunciation Dictionary	Defining syllables ( <i>Block 3</i> )
Omnizart Vocal-Contour	Vocal-contour transcription ( <i>Block 4</i> )
nllb-200-distilled-600M	Baseline model for translation ( <i>Block 5</i> )
nllb-200-distilled-600M (fine-tuned)	Fine-tuned model for translation ( <i>Block 5</i> )
pyKAKASI	Convert kanji to hiragana ( <i>Block 6</i> )
Nagisa	Word separation ( <i>Block 6</i> )
Synthesizer V	SVS ( <i>Block 7</i> )

- No standard evaluation metric for SV2SVT or song translation exists.
  - Mean opinion score test of 6 questions with 6 native Japanese speakers, testing both the baseline and fine-tuned translation models. Examples used in evaluation can be found at [silasantonisen.github.io/polysinger](https://silasantonisen.github.io/polysinger).
- 1 = *Very Poor*, 2 = *Poor*, 3 = *Neutral*, 4 = *Good*, and 5 = *Very Good*.

ID	Question	Baseline	Fine-tuned
Q1	How much sense do the lyrics make?	2.53 ± 0.49	2.17 ± 0.46
Q2	How natural is the Japanese used in the lyrics?	2.57 ± 0.48	2.30 ± 0.48
Q3	How well is the meaning of the original lyrics preserved?	2.47 ± 0.44	2.10 ± 0.44
Q4	How singable are the generated lyrics?	2.40 ± 0.41	2.23 ± 0.44
Q5	How well are the lyrics and melody aligned?	2.50 ± 0.52	2.10 ± 0.40
Q6	What is the overall quality of the generated Japanese singing?	2.33 ± 0.45	2.13 ± 0.41

- High variance in opinion scores.
- No statistically significant differences between the baseline and fine-tuned translation systems.
- Most recurring comments from participants revolved around quality and naturalness of the Japanese lyrics, e.g., incorrect pronunciations or word separations, not necessarily the translation.

## CONTACT INFORMATION

E-mail: [santon@ugr.es](mailto:santon@ugr.es)  
Signal Processing, Multimedia Transmission and  
Speech/Audio Technologies Group, SigMAT.  
Dept. of Signal Theory, Telematics and Communi-  
cations, University of Granada, Spain.