# E2E-S2S-VC:
# End-to-end sequence-to-sequence voice conversion

*Takuma Okamoto[1], Tomoki Toda[2,1], and Hisashi Kawai[1]*

[1]National Institute of Information and Communications Technology, Japan, [2]Nagoya University, Japan

## Demo samples and source code

Demo samples: Hi-Fi-CAPTAIN corpus for Japanese used in experiments

Source code based on ESPnet2-TTS
- Recipe for CMU-ARCTIC corpus
- Recipe for Hi-Fi-CAPTAIN corpus used in experiments

https://ast-astrec.nict.go.jp/demo_samples/e2e-s2s-vc/

## Hi-Fi-CAPTAIN: Released!

High-fidelity and high-capacity conversational speech synthesis corpus developed by NICT

1 female and 1 male (English): 14K utts (parallel: 13K)
1 female and 1 male (Japanese): 19K utts (parallel: 18.5K)
ESPnet2-TTS recipe for JETS-based E2E-TTS

https://ast-astrec.nict.go.jp/en/release/hi-fi-captain/

## 1. Introduction

- **Voice conversion (VC) methods**
  - Framewise VC based on frame-by frame conversion
    - Parallel data not required
    - Difficult to convert duration and prosody between source and target speakers
    - End-to-end models have been investigated (e.g. NVC-Net)
  - Sequence-to-sequence (S2S) VC
    - Parallel data required
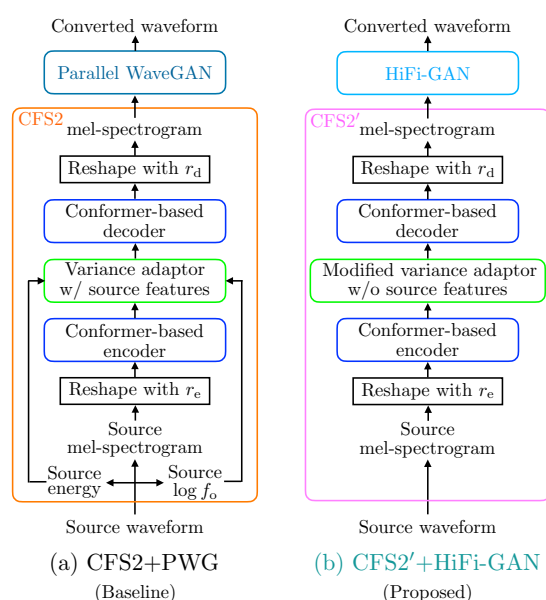    - Can convert duration and prosody by S2S manner
- **Baseline: non-autoregressive (AR) S2S-VC: CFS2+PWG**
  - Features
    - Conformer-Fastspeech 2 (CFS2)-based non-AR conversion model with Parallel WaveGAN (PWG) neural vocoder
    - Faster and more stable by non-AR structure compared with conventional Transformer-based AR models
  - Four problems
    1. Three models (teacher Transformer, CFS2, PWG) are separately trained -> they cannot be jointly optimized
    2. Unstable alignment due to teacher AR Transformer
    3. HiFi-GAN neural vocoder outperforms PWG
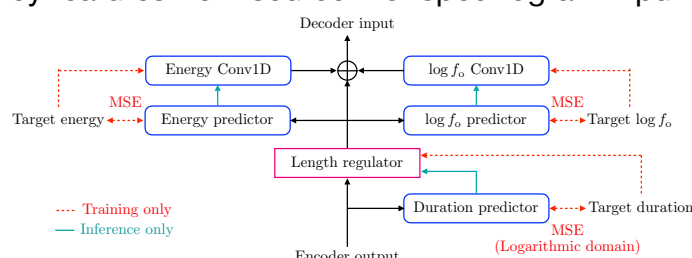    4. Energy and fundamental frequency of source speaker required

## 2. Extended model

- **CFS2'+HiFi-GAN**
  - CFS2': CFS2 with modified variance adapter without source energy and fundamental frequency features

(a) CFS2+PWG (Baseline)

(b) CFS2'+HiFi-GAN (Proposed)

- Modified variance adapter predicts target energy and fundamental frequency features from source mel-spectrogram input
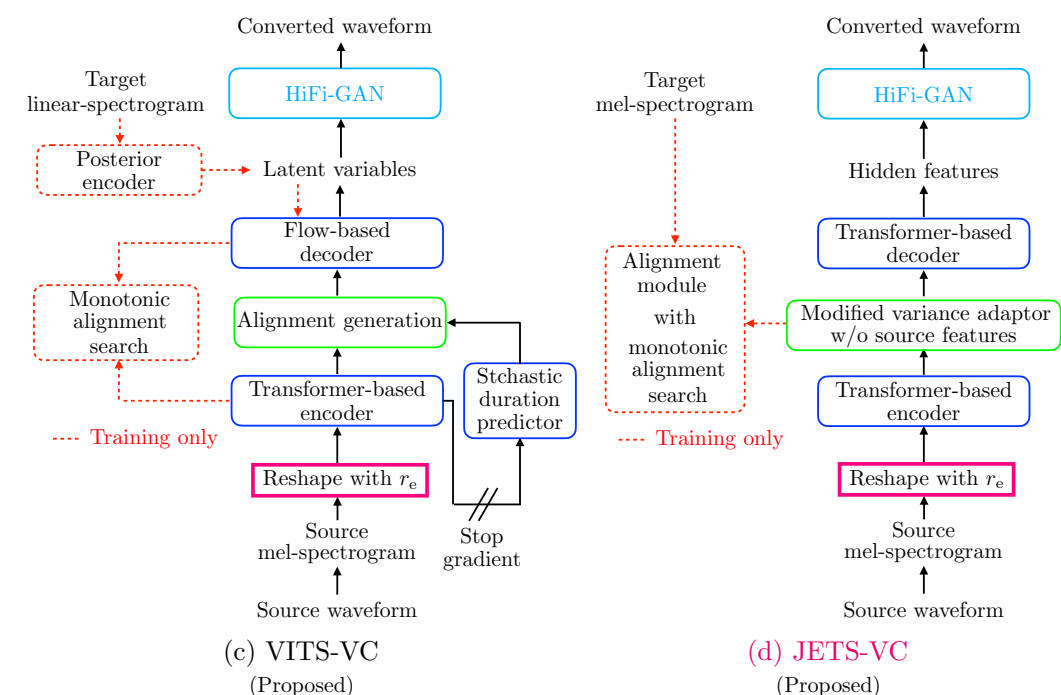
## 3. Proposed methods

- **End-to-end text-to-speech (E2E-TTS) models: VITS and JETS**
  - VITS: VAE + flow + HiFi-GAN + monotonic alignment search (MAS)
  - JETS: Fastspeech 2 + HiFi-GAN + MAS
- **Proposed E2E-S2S-VC: VITS-VC and JETS-VC**
  - Introducing VITS and JETS for E2E-TTS models into S2S-VC
    - Source mel-spectrogram input can be directly converted to target speech waveform with a single neural network
    - Using a reduction factor only for encoder to successfully train MAS for VC
    - Can solve all the four problems in baseline model

(a) CFS2+PWG (Baseline)

(b) CFS2'+HiFi-GAN (Proposed)

(c) VITS-VC (Proposed)

(d) JETS-VC (Proposed)

End-to-end models

## 4. Experiments

- **Experimental conditions**
  - Dataset: Parallel 1,000 utterances for Japanese in Hi-Fi-CAPTAIN
    - Training: 950 utts, Validation: 25 utts, Evaluation: 25 utts
  - Sampling frequency: 24 kHz
  - Objective evaluation criteria: MCD, log$f_o$RMSE, CER and RTF
  - Subjective evaluation criteria (N=20): MOS, speaker similarity
- **Results of experiments**

| Method | Male $\longrightarrow$ Female | | | Female $\longrightarrow$ Male | | | |
|---|---|---|---|---|---|---|---|
| | MCD [dB] | log $f_o$ RMSE | CER [%] | MCD [dB] | log $f_o$ RMSE | CER [%] | RTF |
| Original | N/A | N/A | 1.0 | N/A | N/A | 1.2 | |
| Baseline: CFS2+PWG | 5.83 ± 0.52 | 0.25 ± 0.07 | 3.4 | 4.74 ± 0.26 | 0.20 ± 0.04 | 4.4 | 3.44 |
| CFS2'+PWG | 5.50 ± 0.45 | 0.24 ± 0.08 | 3.0 | 4.76 ± 0.23 | 0.18 ± 0.06 | 6.8 | 3.41 |

### Naturalness

| | Male → Female | Female → Male | Average |
|---|---|---|---|
| Original | 4.79 | 4.39 | 4.59 |
| CFS2+PWG | 3.75 | 3.89 | 3.82 |
| CFS2'+HiFi-GAN (ft) | 3.95 | 3.56 | 3.75 |
| $r_e = 2$ | 4.27 | 3.77 | 4.02 |
| JETS-VC $r_e = 3$ | 4.05 | 4.05 | 4.05 |

| Similarity [%] | 68±7 | 59±7 | 73±6 | 68±7 | 67±7 | 80±6 | 69±6 | 79±6 | 67±5 | 69±5 | 71±4 | 73±4 |