



Experiment settings

dataset	Voice Actre 100 utteran
# of train and test	
sampling frequency	
speech feature	F0 / Pow
DNN	Three-layer Perceptro
CPU	Intel (R) Coi

<u>Computation cost for the proposed VC</u> RTF (Real-Time Factor): a value of the computation time divided by the length of the input waveform

	Analysis	Conversion	Synthesis
Processing Time [ms]	0.85 ± 0.092	1.12 ± 0.51	0.57 ± 0.10
RTF	0.17	0.22	0.11

 \rightarrow The total RTF was <u>smaller than 1.0</u>, so the proposed VC system operates in real time.

ress Corpus [y_benjo et al, 2017] ices, ~12 min, two male speakers

90 / 10 utterances

16 kHz

ver/ MCEP (up to 40 dim) / BAP

n, Input-Output Residual [Saito et al., 2017]

re (TM) i7-3770K CPU @ 3.50 GHz





35 listeners participated in evaluation on 10 utterances recorded at a noisy place. (preference AB / XAB)

 \rightarrow The physical mask is effective for naturalness of speech under noisy environment.

We'll present a further application of this mask in ACM UIST 2019 posters. preprint: <u>https://rikky0611.github.io/resource/transvoice_uist2019poster_paper.pdf</u>

5	S

		with Mask-sh	laped Device	without Ma	sk-shaped Device		
ock							
	Naturalness	33.4%		66.6%			
ine				CO (-0)			
eech	Individuality	30.6%		69.4%			
2.	09	% 25	% 50)% 75	5% 100%		
	Preference Score						
	<u>train on noisy data</u>						
		with Mask-sł	naped Device	without Ma	sk-shaped Device		
	Naturalness	55	.7%	44	4.3%		
				_			
(Individuality	41	.7%	58	8.3%		