

1 **Response to reviews of paper 3681 – “Blow: a single-scale hyperconditioned flow for non-**  
2 **parallel raw-audio voice conversion”.**

3 We are very grateful to all reviewers for their positive assessment of our work and their constructive feedback. We have  
4 considered all minor remarks and typos and already updated our version of the manuscript accordingly.

5 Two reviewers comment on the forward-backward conversion procedure and our design hypothesis that the  $\mathbf{z}$  space  
6 should be condition-free (or identity-neutral in the case of voice conversion). We believe the suggestions made in this  
7 regard are very interesting, and thank the reviewers for that. We have acted in two ways to improve our work. On  
8 the one hand, we have rephrased the “Forward-backward conversion” section (Sec. 4.3), removing some terms that  
9 complicated the understanding of the procedure and improving the wording of the remaining ones. On the other hand,  
10 following the reviewers’ suggestion, we have performed an additional experiment to bring evidence that the  $\mathbf{z}$  space  
11 is identity-neutral (the details of this experiment are already included in our current version of the manuscript). Our  
12 reasoning for the experiment that, if  $\mathbf{z}$  vectors carried information about the speaker, then a powerful-enough classifier  
13 on top of  $\mathbf{z}$  vectors should be able to perform speaker identification much better than random chance, which is 1.1% for  
14 the considered data set and split. In the extreme case, if  $\mathbf{z}$  vectors carried as much information about speaker identity as  
15 the original  $\mathbf{x}$  vectors, the classifier should approach the accuracy of the classical-feature linear classifier we use for the  
16 Spoofing metric, which is 99.3% for the considered data set and split. Hence, on top of  $\mathbf{z}$  vectors, we trained both a  
17 random forest classifier with 50 estimators and a multilayer perceptron with 3 layers, and obtained accuracies of 1.8  
18 and 1.4%, respectively; both marginally above random chance. This gives us an indication that that the  $\mathbf{z}$  space has  
19 limited speaker identity information or that, in the best of cases, such information is not apparent or complex to extract.

20 Finally, one reviewer points out that we should compare to “AUTOVC: Zero-shot Voice Style Transfer with Only  
21 Autoencoder Loss (ICML 2019)”. We thank the reviewer for pointing us out to this very recent and interesting work,  
22 and we will include it in our “Related work” section (Sec. 2). At the moment of writing this author feedback response,  
23 we are working with the provided code to make it work under our setting in order to, in the future, be able to provide  
24 additional results.