

1 **General:** Thank you for your feedback. The main motivation/intuition behind our approach is the excellent performance  
2 we achieved with U-Nets for image segmentation. The theoretical/conceptual contribution is to exploit the similarity  
3 between sleep staging – and time series segmentation in general – and image segmentation. Still, U-Time is not simply  
4 a 1D U-Net (e.g., because of the *segment classifier*, dilated convolutions, normalization, NN-upsampling, ...).

5 We did not run our DeepSleepNet implementation as a baseline on all datasets, because of lacking compute resources  
6 and, given the extensive experiments reported in the supplement, because we were (and still are) sure that U-Time  
7 would clearly outperform it. We now understand the value of explicitly stating the additional baseline results and are  
8 currently running the corresponding experiments. We would like to stress that we are comparing a system developed on  
9 a single design dataset (S-EDF-39) across several datasets. In contrast, the baselines come from papers that propose one  
10 system for one dataset. Thus, the cited results are prone to (unintentional) method overfitting. This is why we regard  
11 our results as a significant advance to automated sleep staging research and clinical practice.

12 The N1 class is hard because it is rare and difficult to clearly separate from other classes, even for human experts as can  
13 be seen in Table 1. N1 lies between wake and deeper non-REM sleep, and the transitions are gradual.

14 **Rev. 1:** Yes, we used only one channel for a fair comparison to existing literature. We are currently evaluating the use  
15 of multiple channels and would be happy to report results in the supplement. We will include a supplementary figure  
16 showing an exemplary confidence score (softmax) output for a  $C = 3$  multi-channel input, see Figure 1 below and our  
17 response to Rev. 2.

18 Details of the layers can be found in the supplement Table S.2. We will add another, larger figure illustrating the  
19 U-shape and the three parts of the architecture. We will make the code available for maximum reproducibility. The  
20 classes in our application are already very imbalanced, so U-Time may also work for very rare events.

21 **Rev. 2:** We are happy to restructure the manuscript according to your suggestions. Regarding confidence scores,  
22 qualitative analysis, and interpretability: It is particularly interesting to inspect the outputs of U-Time when the (freely  
23 adjustable) segmentation frequency is set to match to the input signal frequency, see Figure 1 below. The sleep stage  
24 scores indeed show human interpretable patterns even on short timescales. We believe that this special property of  
25 U-Time will allow for a better analysis of sleep stage transitions in healthy and diseased populations. We will discuss  
26 this in the main text including cases where the model fails to predict the true sleep stages and add the figure below to  
27 the supplementary material.

28 AASM stands for *American Academy of Sleep Medicine* (see line 129, sorry, we forgot to add the abbreviation).  $T$   
29 is the number of fixed-length connected segments (each typically 30s) input to the model (line 68), which we will recall in  
30 line 169.  $B$  is the batch size, which we will introduce properly.

31 **Rev. 3:** Regarding architecture choice: The U-Net part of  
32 our architecture is based on an architecture that worked  
33 extremely well in medical image segmentation across a  
34 wide range of problems. There is a paper at the upcom-  
35 ing MICCAI showing this, which we cannot cite without  
36 revealing authors of the current submission (we can share  
37 a preprint with the area chair). In our submission (lines  
38 91–97), we discuss the sizes of the kernels in relation to  
39 the time windows the layers see – an important design  
40 criterion. We selected them based on our physiological  
41 understanding of sleep staging. We regard it as a big  
42 advantage that we did not extensively tune our architec-  
43 ture to the tasks. The fact that this was not necessary  
44 demonstrates the soundness of the basic approach and the  
45 robustness of the implementation. It is important for us  
46 that our results are not artefacts resulting from (uninten-  
47 tional) overfitting through architecture / hyperparameter  
48 tuning. In contrast, as shown in the supplement, we un-  
49 successfully tried to tune competing methods to reach the performance of U-Time. Tuning the U-Time architecture and  
50 hyperparameters may improve the results. However, we assume that adding different channels or other input modalities  
51 is more important.

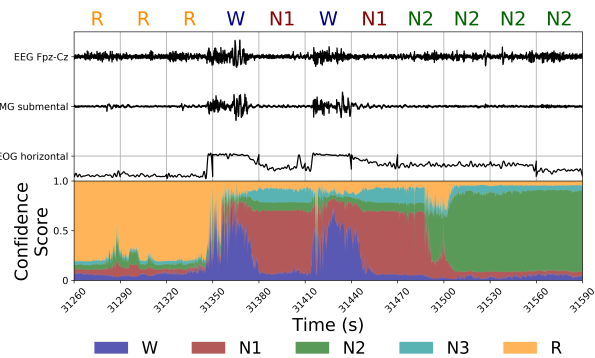


Figure 1: U-Time confidence scores (softmax output) over  $T = 11$  segments (30s each) for  $C = 3$  input channels (EEG, EMG and EOG). The freely adjustable segmentation frequency is set to match the input signal frequency.

52 Confusion matrices (CMs) for U-Time on all datasets are given in the supplement, and we are happy to add CMs for  
53 other methods. We will improve the colour palette as suggested. See the new figure above for an illustration of the  
54 confidence scores (which we will explain better in the main text). We used 5-fold cross-validation on the design dataset  
55 to choose between two loss functions (cross entropy and the dice loss).