

# Improving AHI Scoring Accuracy Using an AI Model for Sleep State Classification from Home Sleep Apnea Testing

H. Zagross Zahawi, S.Æ. Jónsson, E. Arnardóttir, E. Finnsson, E. Erlingsson, K. Montazeri, P.B. Sigmarsdóttir, H.D. Hlynsson, M.E. McPhee Christensen, J.S. Ágústsson

Nox Research, Nox Medical ehf, Reykjavík, Iceland

Presented at Nordic Sleep Conference 2023, Iceland

**nox**  
RESEARCH

thora@noxmedical.com

## Introduction

Home-sleep-apnea-testing (HSAT) devices typically provide limited information on sleep stages, which can limit the accuracy of Apnea-Hypopnea Index (AHI) scoring, a crucial metric used in diagnosing sleep apnea. We present an artificial intelligence (AI) model capable of classifying sleep stages (Wake, REM, NREM) and estimating total sleep time (TST) from HSAT data, that can improve AHI scoring accuracy without needing full polysomnography (PSG). The model uses temporal information from respiratory inductance plethysmography (RIP) and activity signals to predict sleep stages.



## Methods

The temporal convolutional network [1] model was developed using RIP and activity data from sleep studies that were manually scored by sleep technicians in numerous sleep centers across five countries. The recordings were divided into training, validation, and testing subsets during model development.

To further evaluate the HSAT model's performance and generalizability, we conducted a final validation on an external dataset with a full PSG setup originating from a separate sleep clinic that had not been used previously. This diverse dataset included sleep studies on 996 adult participants with: no, mild, moderate, and severe sleep apnea. The dataset served to assess the model's ability to handle new and unseen data, which is integral for real-world applications.

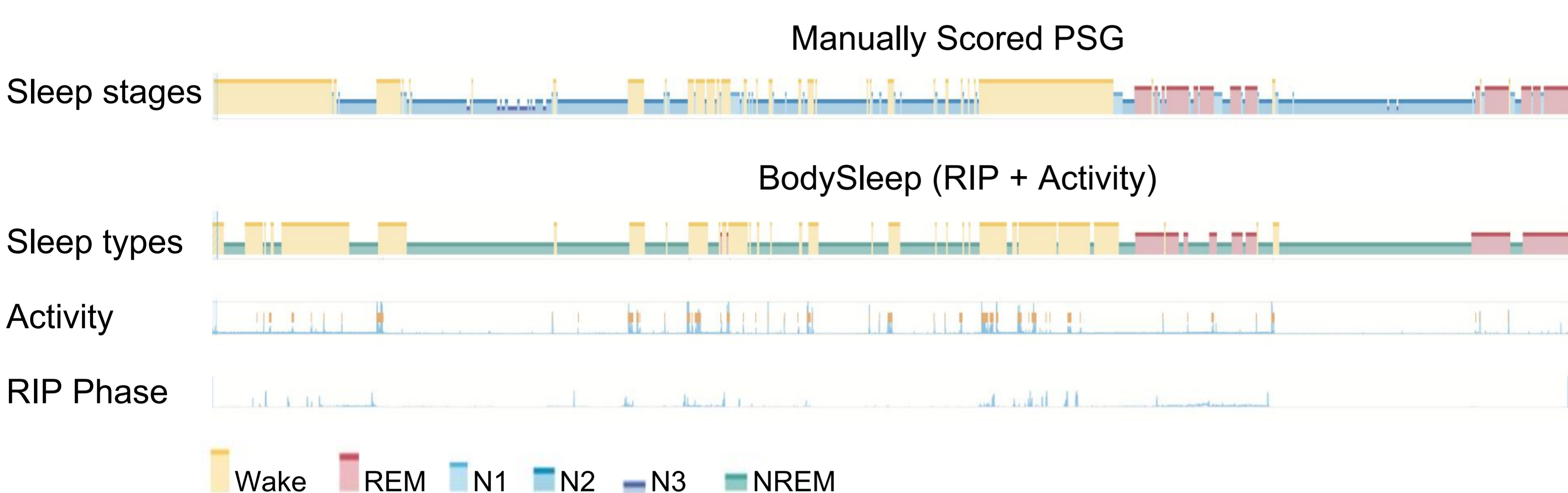
As the dataset used for validation contained full PSG data, it was further possible to compare the accuracy of sleep staging in the present AI model using HSAT (RIP and activity) signals only, to the accuracy of automated full PSG sleep staging (which also includes EEG/EOG/EMG signals). Agreement of the HSAT model with manual scoring for AHI classifications was also investigated.

## Results

The AI model of HSAT signals was compared to manual scoring of sleep stages in 996 sleep studies containing 825,535 valid sleep epochs. Sensitivity, specificity and accuracy of sleep staging was calculated.

**Table 1: Agreement in sleep stages between the HSAT model and manual scoring with accuracy of a full PSG AI model for comparison**

	Sensitivity	Specificity	Overall Accuracy
<b>Wake</b>	<b>80%</b>	<b>95%</b>	<b>92%</b>
Full PSG Model*	64% (53 - 74%)	97% (94 - 98%)	94% (92 - 96%)
<b>REM</b>	<b>81%</b>	<b>98%</b>	<b>96%</b>
Full PSG Model*	79% (66 - 87%)	97% (93 - 98%)	94% (91 - 96%)
<b>NREM</b>	<b>92%</b>	<b>82%</b>	<b>89%</b>
Full PSG Model*	93% (88 - 96%)	80% (72 - 86%)	90% (86 - 92%)

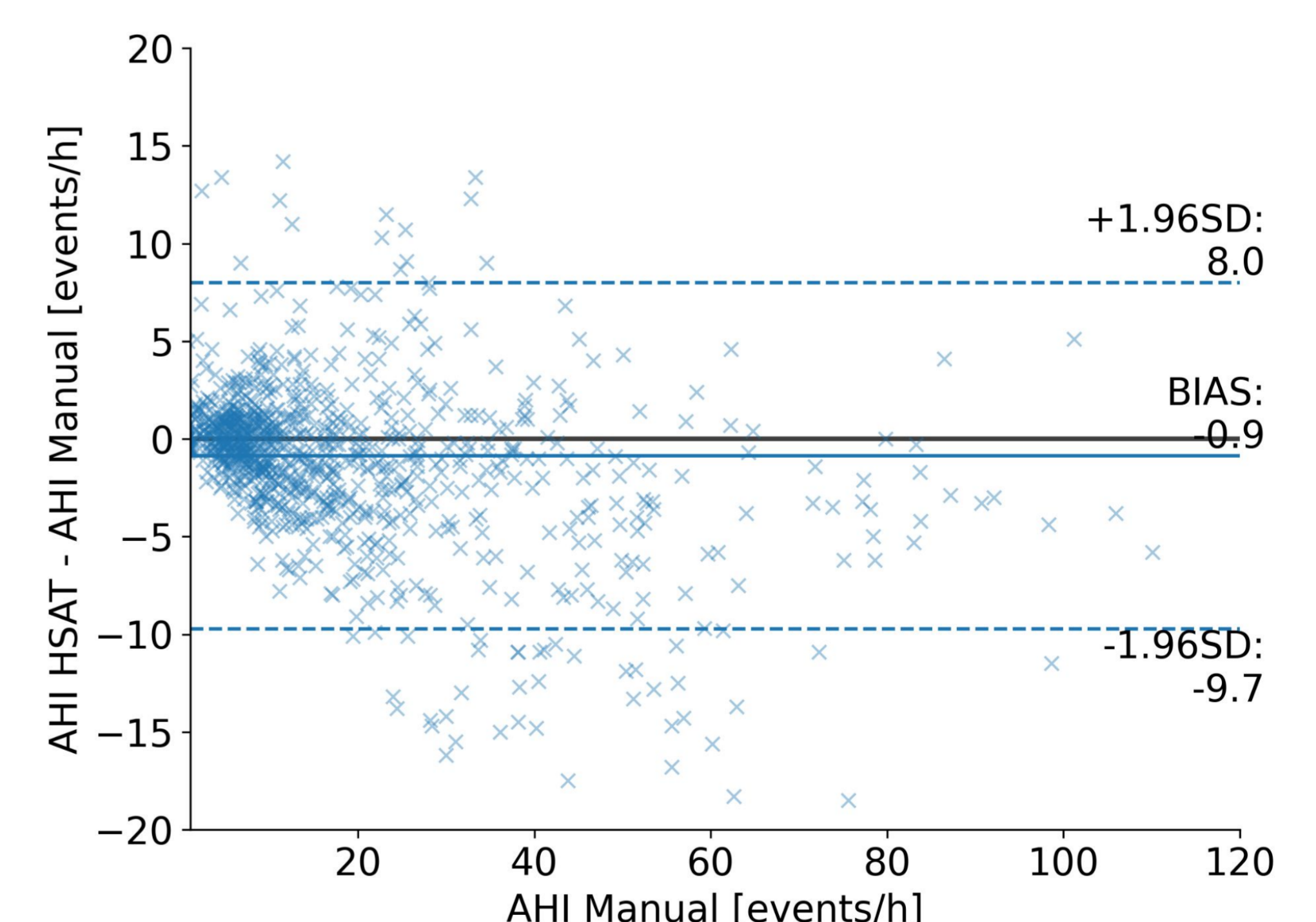


**Figure 1: An example of the sleep profile in a patient with sleep apnea. The figure shows a single sleep study scored manually and scored automatically using RIP and activity signals to allow for comparison of these two methods in capturing Wake, REM, and NREM periods.**

Based on the AI model, patients were classified by AHI for which sensitivity, specificity and accuracy were calculated with manual scoring as the reference standard.

**Table 2: Agreement of AHI classifications from the HSAT model with manual scoring (N = 996 Patients)**

	Sensitivity	Specificity	Overall Accuracy
<b>AHI ≥ 5</b>	<b>96%</b>	<b>92%</b>	<b>95%</b>
<b>AHI ≥ 15</b>	<b>88%</b>	<b>96%</b>	<b>93%</b>



**Figure 2: Bland-Altman for the HSAT model compared to manual scored AHI**

## Conclusions

Automatic scoring of sleep stages from HSAT (breathing and activity) signals shows similar accuracy to a full PSG AI model in determining sleep stages, when compared to manual scoring. The present HSAT model also showed accurate AHI classification, suggesting that it can accurately and quickly score clinically important sleep metrics.

1. Bai, S., Kolter, J. Z., & Koltun, V. (2018). An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*.

\* Noxturnal 6.1