

Motivation: Current AI models used for sleep stage classification remain largely opaque black boxes that healthcare providers cannot trust or validate. These systems are potentially vulnerable to adversarial attacks that could compromise patient safety or be used to manipulate health data.

① Utilize DREAMT dataset

(Wang, Ke, et al.)³



Technician-annotated sleep stage labels



Wrist-worn wearable signals

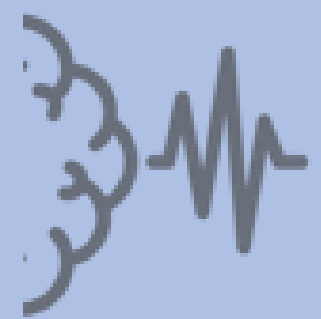
② Pre-process data



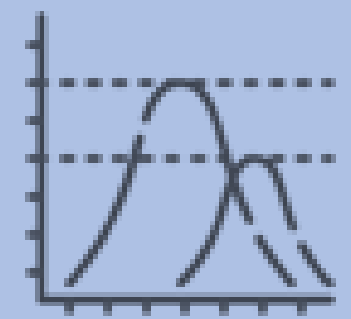
Divide data into 30 second epochs



Statistical feature extraction



Signal processing techniques



Gaussian Filters

③ Structure input data

Pre-extracted features

Raw data

④ Select models

Support Vector Machine

(Logacjov, Aleksej, et al.)⁴

Baseline for evaluating robustness and explanation stability

LightGBM + Long Short-Term Memory

(Wang, Ke, et al.)³

Captures population-level patterns and subject-specific temporal dependencies

Convolutional Neural Network + Recurrent Neural Network

(Korkalainen, H, et al.)⁵

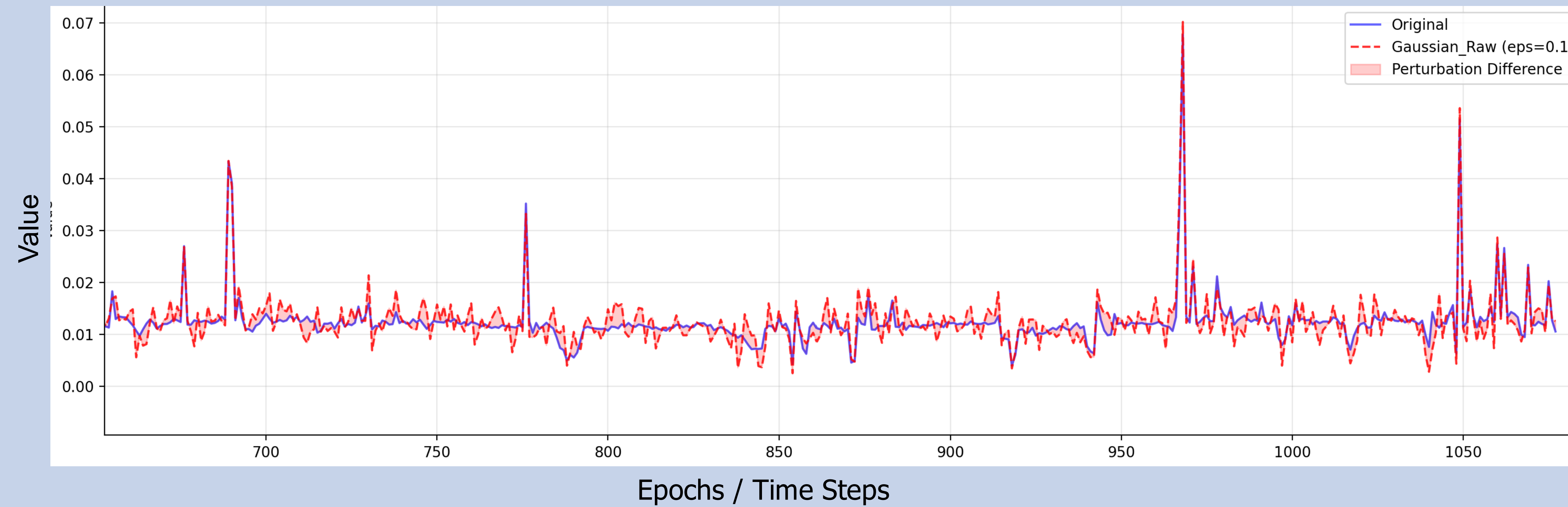
Models complex temporal and cross-sensor feature interactions

5a) Apply adversarial attacking techniques on all models to test robustness

F1 Scores on Selected Models and Attacking Techniques

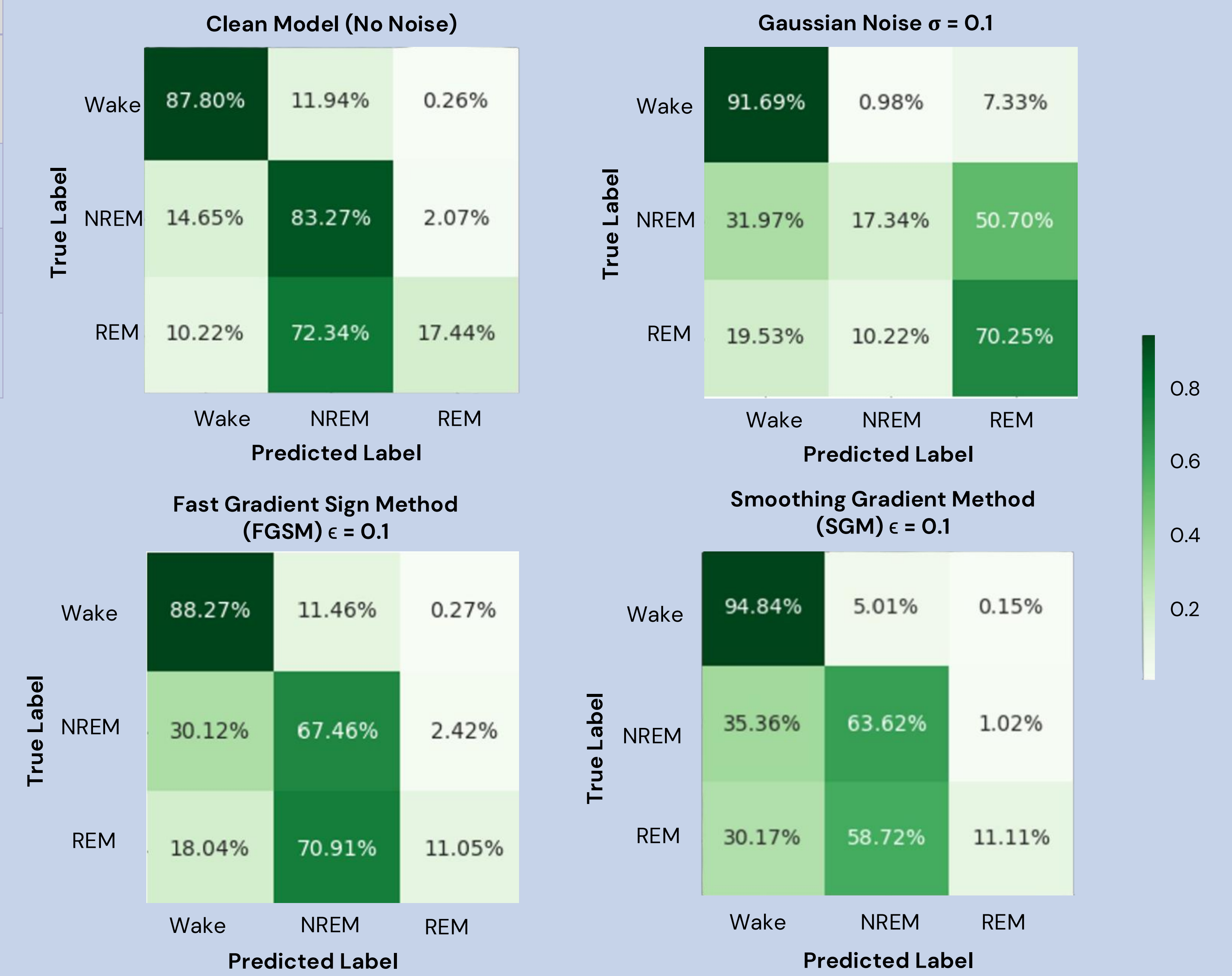
	Baseline	Gaussian Noise ($\sigma = 0.05$)	Gaussian Noise ($\sigma = 0.1$)	Gaussian Noise ($\sigma = 0.2$)	FGSM ($\epsilon = 0.1$)	Smoothing ($\epsilon = 0.1$)
SVM (2 Stage)	0.81	0.81	0.81	0.80	0.75	-
LightGBM + biLSTM (3 Stage)	0.64	0.47	0.45	0.44	0.55	0.55
CNN + RNN (3 Stage)	0.38	0.36	0.36	0.36	-	-

Accelerometry Signal with Gaussian Noise Injection $\sigma = 0.1$, Participant 078



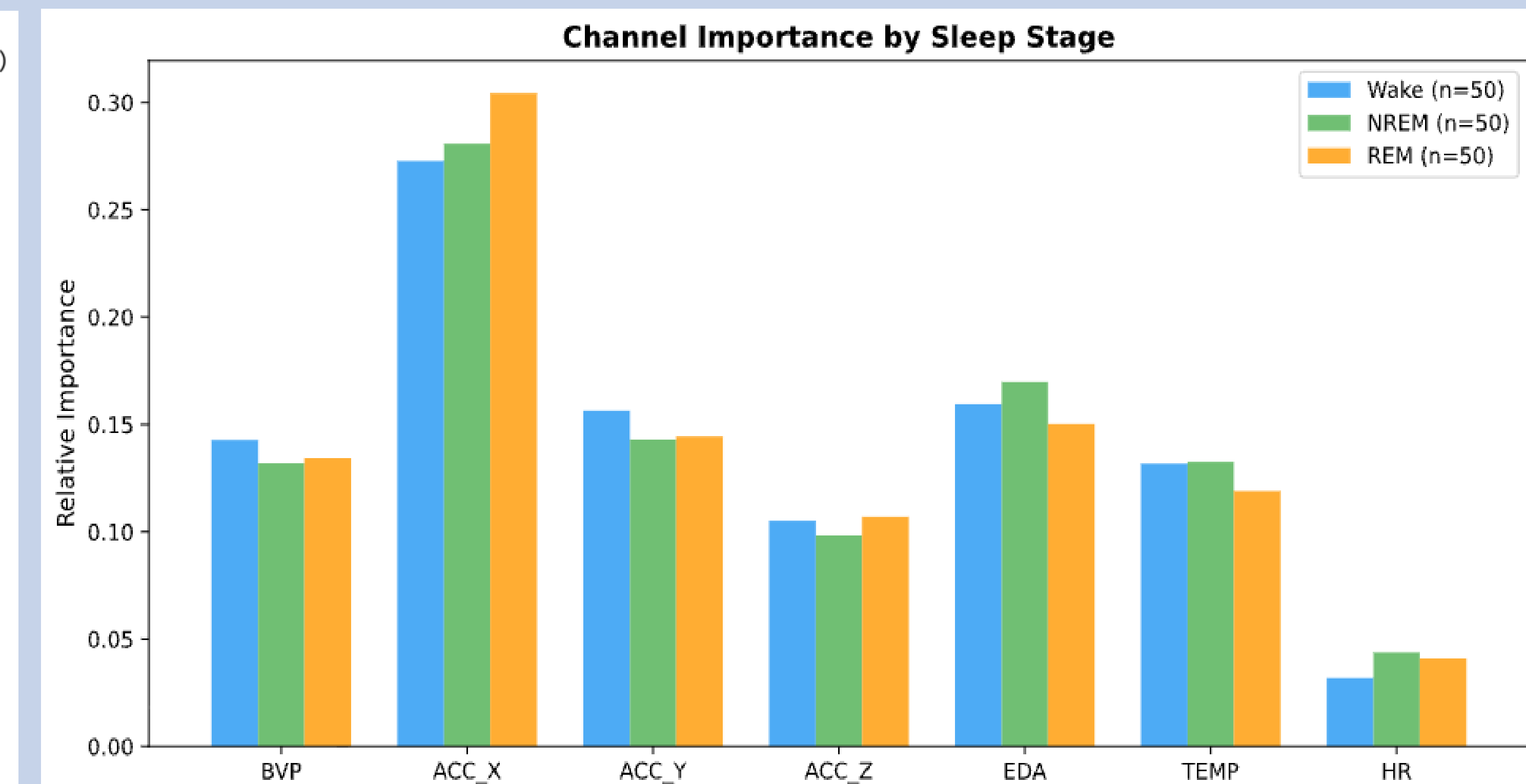
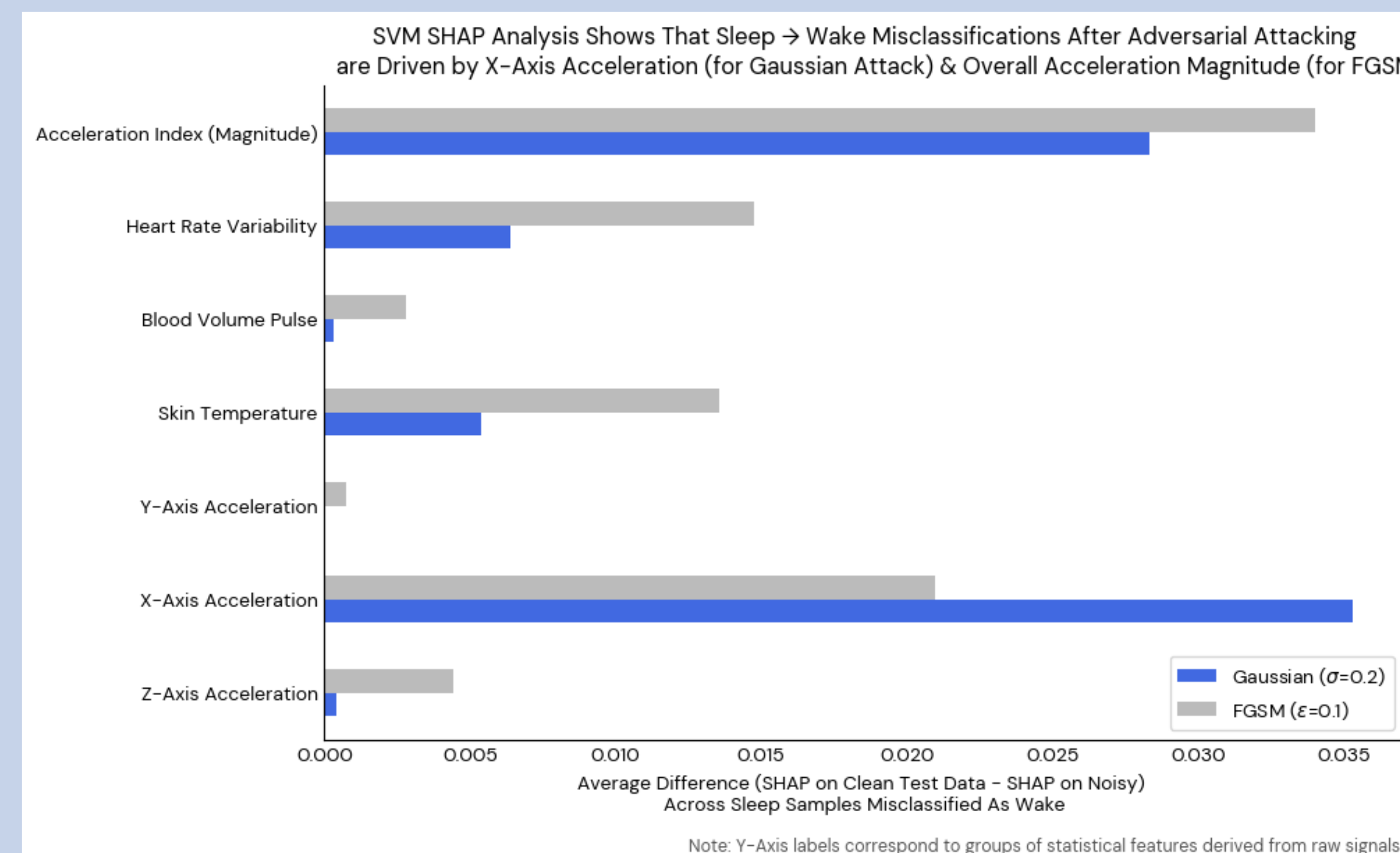
LightGBM - LSTM Sleep Stage Classification

Injecting noise worsens LightGBM - LSTM model performance overall. Gaussian noise on raw signal data pushes model predictions away from NREM to REM while gradient based methods on the LSTM push model towards more Wake predictions.



5b) Use explainability techniques to understand how feature attributions change with noise injection

Differencing SHAP values for samples that are misclassified by the SVM model after an adversarial attack helps us see what feature groups push less towards a sleep classification (which is the true classification) after the attack, and thus what groups are most impacted by noise and drive misclassification.



Grad-CAM was used on the CNN-RNN sleep staging model to evaluate feature importance by channel/signal. Results were averaged across 50 balanced samples per class. X-Axis acceleration (ACC_X) emerges as the dominant feature across all three sleep stages. Its relative importance is consistent with the two other models' results.

³ Wang, Ke, et al. "DREAMT: Dataset for Real-time sleep stage Estimation using Multi-sensor wearable Technology" (version 2.1.0). PhysioNet (2025). RRID:SCR_007345. <https://doi.org/10.13026/7r9r-7r24>
⁴ Logacjov, Aleksej et al. "A Machine Learning Model for Predicting Sleep and Wakefulness Based on Accelerometry, Skin Temperature and Contextual Information." Nature and science of sleep vol. 16 699-710. 6 Jun. 2024. doi:10.2147/NSS.S452799
⁵ Korkalainen, Henri, et al. "Deep Learning Enables Sleep Staging from Photoplethysmogram for Patients with Suspected Sleep Apnea." Sleep, vol. 43, no. 11, 21 May 2020. <https://doi.org/10.1093/sleep/zaa098>.