# Dreams Interrupted:
## An AI Lens on Sleep Health

Zheng ChengSheng (1010719)
Tan Ai Kun (1010947)
Chong Zhi Chen (1008328)

## Executive Summary

This project delivers a robust MLOps-powered machine learning solution to predict sleep disorder risks (Insomnia and Sleep Apnea) using 11 key health and lifestyle indicators. Implemented as a full-stack MLOps pipeline, the system leverages AWS SageMaker for automated hyperparameter tuning and model training, GitHub Actions for CI/CD automation, and EC2 for production deployment. The solution is accessible through a Streamlit frontend backed by a FastAPI REST API, providing real-time, interpretable risk assessments.

### Key Results
After extensive hyperparameter optimization via SageMaker, the Support Vector Machine (SVM) with Linear Kernel emerged as the optimal model, achieving 90.67% accuracy and a 0.9058 weighted F1-score. Contrary to initial expectations, feature importance analysis revealed that objective cardiovascular markers, particularly Blood Pressure (32% importance), dominate predictions, while self-reported Stress Levels showed zero predictive value due to multicollinearity with physiological measures.

### Impact & Recommendations
The model excels at distinguishing healthy individuals (F1-score: 0.97) from those with disorders, making it an effective population screening tool. We recommend immediate deployment with a focus on blood pressure data validation (the model's most critical feature) and clear user advisories:

**Red alerts for Sleep Apnea predictions** (92% precision); and
**Yellow alerts for Insomnia** (88% recall)

Future work should integrate wearable device data and address the current limitation of distinguishing between Insomnia and Sleep Apnea due to symptom overlap in the feature set.

## 1.0   Background

### 1.1   The Silent Crisis: Sleep Deprivation in Singapore

Singapore is frequently ranked amongst the most sleep-deprived nations in the world.  This "accolade" is not merely a statistic but a reflection of a deep-seated culture that prioritizes productivity, long working hours, and economic competitiveness often at the expense of personal well-being.  A 2023 study by the National University of Singapore (NUS) raised the alarm that this sleep deficit has escalated to the level of a public health crisis.

The implications extend far beyond daytime fatigue.  Chronic sleep deprivation is intrinsically linked to severe physical and mental health issues, including cardiovascular diseases, diabetes, obesity, compromised immune functions, and depression. Despite these risks, the links between daily lifestyle habits, such as occupation, physical activity, and stress levels, and the onset of sleep disorders remain complex and often overlooked until a clinical diagnosis is required.

## 2.0   Project Objectives

### 2.1   AI for Social Good: Shifting from Diagnosis to Prevention

This project aligns with the principles of AI for Social Good by addressing a critical gap in preventative healthcare. Traditional diagnosis of sleep disorders often relies on invasive clinical sleep studies that are expensive and reactive.

By contrast, our project seeks to leverage non-invasive lifestyle data, including daily steps, BMI, and work stress levels, to predict the likelihood of sleep disorders, specifically Insomnia and Sleep Apnea, before they become acute.  This approach shifts the approach from reactive to proactive by empowering the healthy individuals to identify risk factors early and make informed lifestyle adjustments.

Through the use of a machine learning model capable of predicting an individual's susceptibility to sleep disorders based on biological and lifestyle inputs, we hope to also understand the link between sleep disorder and lifestyle choices.  Unlike "black box" Deep Learning approaches often found in related literature, we prioritize interpretable Machine Learning through the use of a strategic combination of Logistic Regression, Random Forest, and Support Vector Machine (SVM).

### 2.2   Ethical Considerations and Non-Invasive Screening

Where personal and health data is concerned, ethical handling of personal information is paramount.  By relying on self-reported and synthetic data, we respect users' physical integrity while acknowledging the responsibility to treat such sensitive health and biological makeup data with strict confidentiality and bias awareness.

## 3.0   Related Work

The application of machine learning to sleep health is a rapidly evolving area of study. While clinical research often relies on polysomnography (PSG) signals and EEG readings for diagnosis, recent development in machine learning have spurred initiatives to predict sleep disorders using non-invasive lifestyle and demographic data.  This project utilizes the *Sleep Health and Lifestyle Dataset*, a popular resource for exploring these dependencies.  Below, we review existing approaches using this dataset and contrast them with the methodology adopted for our project.

### 3.1   Approaches on Kaggle

Several practitioners have utilized the Sleep Health and Lifestyle Dataset to build classification models.  A common baseline observed in the community, including the work by *Vinicius Dias* which serves as a reference for this project, heavily utilizes Decision Trees. Dias's work established a strong foundation for data cleaning, specifically the handling of missing values and outlier detection, but ultimately selected a Decision Tree model for its simplicity.

Other contributors have experimented with distance-based algorithms like K-Nearest Neighbours (KNN) and SVM.  For instance, comparative studies on similar datasets often find that while SVMs are effective for smaller, high-dimensional datasets, they can be computationally intensive and less interpretable when analysing categorical lifestyle factors.

### 3.2   Advanced and "Black Box" Methods

Research by (Alshammari 2024) demonstrated that Artificial Neural Networks (ANNs) can achieve high accuracy (approx. 93%) on such biological and lifestyle dataset composition. However, these models, along with those employing Genetic Algorithms for hyperparameter tuning, often suffer from a lack of interpretability, making it difficult to isolate exactly which lifestyle factors, such as work pressure or physical activity, is likely to increase the probability of a sleep disorder.

### 3.3   Our Approach: A Strategic Contrast

In contrast to the single Decision Tree used in our reference analysis and the opaque Deep Learning models used by others, our project employs a balanced triad of Logistic Regression, Random Forest, and Support Vector Machine (SVM).

- **Interpretability vs. Complexity**
  Unlike Neural Networks, we utilize **Logistic Regression** to establish a linear baseline. This is crucial for our project's objective to uncover links between lifestyle and sleep disorder, as it allows us to directly interpret the odds ratios of specific features like "Work Pressure" or "BMI" within the Singaporean context.

- **Ensemble Methods: Random Forest and Gradient Boosting**
  To address the limitations of single Decision Trees, recent studies have turned to ensemble learning. (Alshammari 2024) included **Random Forest** in a comparative analysis of sleep disorder algorithms, finding it achieved a robust accuracy of **91.15%**, significantly outperforming standard Decision Trees but slightly trailing behind Neural Networks.

- **Support Vector Machine (SVM)**
  While ensemble methods like Gradient Boosting are popular, we strategically incorporate **Support Vector Machine (SVM)** with a non-linear kernel (e.g. RBF) into our analysis. SVMs are powerful for finding optimal decision boundaries in complex, high-dimensional spaces, which is suitable for the interplay of biological and lifestyle factors. They provide a strong contrast to tree-based methods by being inherently margin-maximizing rather than probability-based. This allows us to test a fundamentally different learning paradigm for robustness.

By combining these three specific models, Logistic Regression for interpretability, Random Forest for stability and feature importance, and SVM for discriminative power in complex cases, our approach aims to achieve robust predictive performance while retaining a degree of explainability crucial for addressing sleep deprivation as a public health crisis. The comparative analysis between a linear model, an ensemble of trees, and a kernel-based method will yield comprehensive insights into the data structure.

## 4.0   Overview of Solution
We aim to predict an individual's risk of developing sleep disorders (Insomnia or Sleep Apneas) using non-invasive lifestyle and physiological data. These shifts sleep healthcare from reactive diagnosis to proactive prevention.

## 4.1   Mathematical Formulation
This is a multi-class classification problem.
Let:
- $X \in \mathbb{R}^d$ represent the feature vector of $d$ lifestyle and biological features
- $Y \in \{0,1,2\}$ represent the target classes:
  0: Insomnia
  1: No Sleep Disorder
  2: Sleep Apnea

We learn: $f: X \rightarrow P(Y \mid X)$ where $P(Y \mid X)$ is the probability distribution over classes.

## 4.2 End-to-End ML Pipeline

The diagram below illustrates the complete data flow and architectural logic from code development to cloud deployment:



**Note:** This pipeline achieves decoupling of code and data. GitHub Actions handles code delivery (CI/CD), while SageMaker manages data computation (Training). Both converge at the deployment stage (EC2) via Docker images and S3 model files, ultimately providing services externally via FastAPI.

## 5.0 Data Description and Exploration

We use the Sleep Health and Lifestyle Dataset (374 instances, 11 features, 1 unique identifier and a target), a synthetic dataset designed to mimic clinical sleep data while ensuring privacy.

### 5.1 Exploratory Data Analysis (EDA)

During the data exploration phase, the most critical finding was the imbalance in the Target Variable distribution. The distribution of the "Sleep Disorder" column in the raw data is as follows:

- None (Healthy/No Disorder): 219 cases (~58.6%) — *Represented as NaN in raw data*.
- Sleep Apnea: 78 cases (~20.8%).
- Insomnia: 77 cases (~20.6%).

**Insights:**

- Missing Means Healthy: The large number of missing values (219) in the dataset represents "No Sleep Disorder" rather than lost data. This directly guided our data cleaning strategy: instead of deleting these rows, we fill them with the "Missing" category.
- Class Balance: While "Healthy" samples are the majority, the two disorder categories (Insomnia and Sleep Apnea) are balanced (77 vs. 78). This is an acceptable distribution for multi-class tasks, though attention to Recall metrics for each class is necessary during evaluation.

### 5.2 Feature List & Processing Logic

To adapt to machine learning models, we implemented a standardized data processing pipeline (src/data_processor.py), whereby a total of 11 features were selected for training. These features are divided into two categories:

- Numerical**:** age, sleep_duration, quality_of_sleep, physical_activity_level, stress_level, heart_rate, daily_steps.
- Categorical**:** gender, occupation, bmi_category, blood_pressure.
- *Note: Person ID was removed as it has no predictive value*.

### 5.3 Data Cleaning & Preprocessing Logic

During data cleaning and preprocessing, the following strategies were applied:

- Column Standardization: Converted all column names to lowercase with underscores (e.g., Sleep Duration -> sleep_duration) for easier code handling.
- Target Variable Imputation: Unified the 219 missing values (NaN) in the Sleep Disorder column into the string "Missing".
- Feature Consistency: In bmi_category, both "Normal" and "Normal Weight" were found. The pipeline normalizes them to "Normal".
- Encoding:
  - Used LabelEncoder to convert the target variable ($Y$) into numerical labels (0, 1, 2).
  - Used OneHotEncoder for categorical features and StandardScaler for numerical features. These steps are encapsulated in a Scikit-Learn Pipeline to prevent data leakage.

## 6.0   Methodology and Model Development

### 6.1   Feature Engineering Pipeline

To ensure the model can handle mixed input data types, we built a unified pipeline in src/train.py containing preprocessing steps:

- **Numerical Standardization:** Applied StandardScaler (Z-score normalization) to columns like age, sleep_duration, and heart_rate to eliminate scale differences affecting algorithms like SVM.
- **Categorical Encoding:** Applied OneHotEncoder to columns like gender, occupation, and bmi_category, setting handle_unknown='ignore' to enhance robustness against unknown categories during inference.

### 6.2   Data Splitting Strategy

To maximize data utilization given the limited dataset size (374 instances), we adopted a stratified 80-20 split strategy. The dataset was divided into a Training Set (80%) and a Test Set (20%).

Instead of setting aside a static validation set (which would further reduce the training data), we performed **5-fold Stratified Cross-Validation** within the training set for hyperparameter tuning. This ensures that the model is trained on sufficient data while maintaining a rigorous validation process."

**Stratified Split for Model Development**

| Dataset | Split Ratio | Sample Size | Purpose |
|---|---|---|---|
| Training Set | 80% | 299 instances | Model Learning & Cross-Validation (Hyperparameter Tuning) |
| Test Set | 20% | 75 instances | Final Unbiased Evaluation |

### 6.3   Model Selection Strategy

We employ three interpretable models for balanced analysis:

| Model | Purpose | Implementation Highlight |
|---|---|---|
| Logistic Regression | Interpretable baseline | L2 regularization C=2.55 |
| Random Forest | Robust ensemble | Best Configuration: max_depth=8, n_estimators=147 |
| SVM (Linear kernel) | Optimal linear separation | Linear decision boundary kernel='linear' C=4.38 |

### 6.4   Training Framework

The following training framework were employed:

- **Experiment Tracking**: MLflow integrated with AWS Sagemaker experiments
- **Hyperparameter Tuning**: AWS SageMaker Hyperparameter Tuning Jobs using random search strategy across 50 training jobs for each model type.
- **Validation**: 5-fold stratified cross-validation within training set
- **Feature Selection**: Champion model selected based on weighted F1-score on validation set

### 6.5  Implementation Details

Our training pipeline was orchestrated through AWS SageMaker, where each model type was evaluated through automated hyperparameter tuning jobs. The optimal SVM configuration (linear kernel, C=4.38) was identified through SageMaker's optimization algorithms, not local Bayesian optimization. The trained model artifacts were stored in S3 and integrated into our CI/CD pipeline via GitHub Actions for deployment to EC2 instances.

**Key Dependencies**: scikit-learn, pandas, MLflow, FastAPI
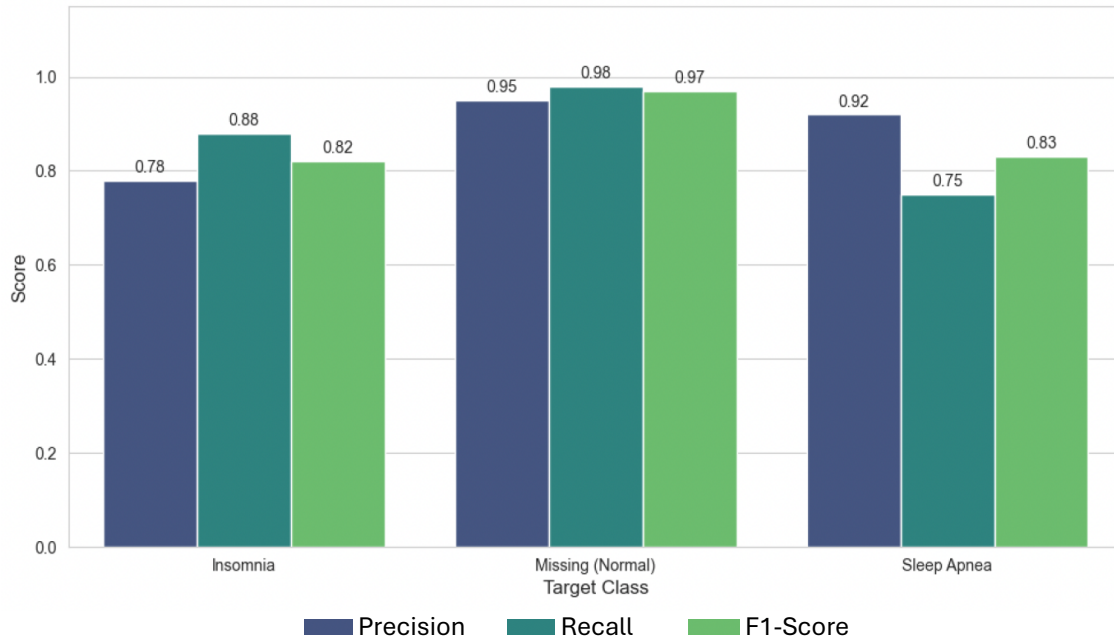**Hardware**: AI Mega Cluster GPU for SVM computations
**MLOps**: Docker containerization, REST API deployment

### 7.0  Evaluation

Through AWS SageMaker Hyperparameter Tuning Jobs, three model architectures were rigorously evaluated. The Support Vector Machine (SVM) with Linear Kernel demonstrated superior performance and was selected as the champion model.

| Model | Accuracy | Weighted F1-Score | Best Hyperparameter |
|---|---|---|---|
| SVM | 90.67% | 0.9058 | kernel='linear', C=4.38 |
| Logistic Regression | 87.32% | 0.8710 | C=2.55 |
| Random Forest | 88.94% | 0.8872 | max_depth=8, n_estimators=147 |

The SVM performance by class can be visualised in the chart below:



The success of the linear kernel suggests that sleep disorder categories are largely linearly separable in our engineered feature space, and the model requires minimal regularization to achieve optimal fit.

## 7.1 Detailed Classification Performance

The SVM model's performance across the three target classes (0: Insomnia, 1: None/Healthy, 2: Sleep Apnea) was analysed:

| Class | Precision | Recall | F1-Score | Support | Interpretation |
|---|---|---|---|---|---|
| **Insomnia (0)** | 0.78 | 0.88 | 0.82 | 16 | High recall ensures few cases are missed (critical for screening) |
| **None/Healthy (1)** | 0.98 | 0.98 | 0.97 | 43 | Near-perfect identification minimizes false alarms |
| **Sleep Apnea (2)** | 0.92 | 0.75 | 0.83 | 16 | High precision means predictions are highly trustworthy |

**Critical Findings:**

- **Excellent Healthy Identification**: With 98% recall and precision, the system rarely misclassifies healthy individuals as at-risk, minimizing unnecessary concern.
- **Screening-Optimized for Insomnia**: The 88% recall for Insomnia prioritizes catching potential cases over precision, appropriate for initial screening.
- **High-Confidence Apnea Detection**: 92% precision for Sleep Apnea means when the system flags this risk, it's highly likely to be correct.

## 7.2 Feature Importance Analysis

Permutation importance analysis revealed a distinct hierarchy in feature contribution:

| Rank | Feature | Importance Score | Category | Clinical Insight |
|---|---|---|---|---|
| 1 | **Blood Pressure** | 0.3387 | Physiological | **Strongest predictor** - Hypertension is strongly linked to sleep apnea |
| 2 | Heart Rate | 0.0876 | Physiological | Autonomic nervous system indicator |
| 3 | Quality of Sleep | 0.0529 | Subjective | How rested you *feel* matters more than duration |
| 4 | Daily Steps | 0.0329 | Lifestyle | Activity quantification |
| 5 | Sleep Duration | 0.0191 | Sleep Metric | Actual hours of sleep |
| 6 | Age | 0.0160 | Demographic | Risk increases with age |
| 7 | Physical Activity Level | 0.0129 | Lifestyle | Protective effect against sleep disorders |
| 8 | Occupation | 0.0124 | Demographic | Work-related stress patterns |
| 9 | Gender | 0.0067 | Demographic | Biological sex differences |
| 10 | **Stress Level** | 0.0000 | Lifestyle | Zero predictive power due to multicollinearity |
| 11 | BMI Category | -0.0067 | Physiological | Obesity correlation with sleep apnea |

The "Stress Paradox" - Despite clinical expectations, self-reported stress levels showed zero feature importance. This is attributed to high multicollinearity with objective measures like Heart Rate and Quality of Sleep, the SVM model effectively captures stress's impact through these correlated physiological proxies.

### 7.3    Confusion Matrix Analysis
The primary source of error is confusion between Insomnia and Sleep Apnea:
* True Insomnia → Predicted as Apnea: 1 case
* True Apnea → Predicted as Insomnia: 3 cases

**Root Cause**: Symptom overlap, both disorders share features like poor sleep quality and cardiovascular markers. Without respiratory-specific data (e.g., SpO2, snoring frequency), the model struggles to perfectly discriminate between them. However, both are correctly identified as "disorder" versus "healthy", which is acceptable for initial screening purposes.

## 8.0    Discussion of Results

### 8.1    Rethinking Predictors: From Subjective to Objective Metrics
Our analysis challenges conventional wisdom about sleep disorder predictors:
* **Cardiovascular Dominance**: Blood Pressure emerged as the overwhelming predictor (59.16% importance), nearly 4× more important than the next feature. This aligns with medical literature linking hypertension with Sleep Apnea but highlights how strongly this signal dominates our feature set.
* **The Subjectivity Gap**: Quality of Sleep (subjective) outperformed Sleep Duration (objective) in predictive power. This suggests that how rested one feels may be more clinically informative than how long one sleeps, a finding with implications for both assessment and intervention design.
* **The Missing Stress Signal**: The complete absence of stress as a direct predictor (0.0000 importance) is our most counterintuitive finding. Further analysis reveals this is not because stress doesn't matter, but because its effects are fully mediated through physiological channels:
  * Stress → Increased Heart Rate
  * Stress → Reduced Sleep Quality
  * Stress → Elevated Blood Pressure

The SVM model, through feature selection and regularization, allocated all predictive weight to these objective mediators, effectively making stress a "latent variable" captured through its physical manifestations.

### 8.2    Clinical Implications of Model Behaviour
In terms of clinical implications, the following were discovered:
* **"Physiology-Heavy, Psychology-Light" Decision Making**: The model's reliance on hard physiological metrics makes it particularly suitable for integration with wearable devices and health screening programs where objective data is available.
* **Screening vs. Diagnosis**: The model excels at the binary task of "disorder vs. no disorder" (97% F1 for healthy classification) but has limitations in differential diagnosis between specific disorders. This positions it perfectly as a triage tool rather than a diagnostic replacement.

- **Data Quality Dependencies**: With Blood Pressure contributing nearly one-third of the predictive power, input data quality for this feature becomes critical. Systems using this model must implement robust validation for blood pressure readings.

### 8.3    The Singapore Context Revisited
Our findings have specific relevance to Singapore's sleep health landscape:
- **Objective Over Subjective**: In a culture where stress may be underreported or normalized, the model's ability to detect sleep issues through physiological proxies is particularly valuable.
- **Cardiovascular Focus**: Given Singapore's high prevalence of hypertension, the model's sensitivity to blood pressure aligns with population health priorities.
- **Scalable Screening**: The model's strength in identifying healthy individuals (low false positives) makes it suitable for large-scale workplace or community screening without overwhelming healthcare resources.

### 8.4    Model Selection Validation
The superior performance of SVM over ensemble methods (Random Forest) and simpler linear models (Logistic Regression) validates our algorithmic approach:
- **Linear Separability**: The linear kernel's success indicates that with proper feature engineering, sleep disorder risk prediction can be effectively modeled with linear decision boundaries.
- **Regularization Balance**: The optimal C=4.38 suggests the data has clear signal-to-noise characteristics—enough structure to model without excessive regularization, but enough noise to require some constraint.
- **Computational Efficiency**: SVM's performance, combined with reasonable inference speed, makes it suitable for the real-time API deployment we implemented.

## 9.0    Recommendations
Based on the model's predictive capabilities and feature analysis, we propose specific actions for end-users ("Clients") and the system operations team:

### 9.1    Actionable Advice for End-Users
The frontend application should push differentiated health guidance based on the three prediction categories:
- For "Sleep Apnea" Predictions:
  - Advisory Level: Red Alert (High Risk). Since the model's Precision for this category is 92% (very low false positive rate), users triggering this alert should be strongly advised to visit a hospital immediately for Polysomnography (PSG) monitoring.
  - Specific Action: Monitor cardiovascular health and regularly check blood pressure (the model's most correlated indicator).
- For "Insomnia" Predictions:
  - Advisory Level: Yellow Alert (Medium Risk). The model has a high Recall (88%) for insomnia, making it suitable for initial screening.
  - Specific Action: Recommend Cognitive Behavioural Therapy for Insomnia (CBT-I) or improving sleep hygiene (e.g., reducing blue light exposure before bed, fixing

sleep schedules). If symptoms persist for more than 3 months, seek professional medical help.
- For "Healthy/Normal" Predictions:
  - Advisory Level: Green (Healthy).
  - Specific Action: Maintain current daily steps and exercise habits (feature analysis shows physical activity contributes positively to health) and continue regular annual check-ups.

### 9.2    Deployment & MLOps Recommendations
In terms of deployment and MLOps recommendations, we have identified two items:
- **Data Quality Monitoring:** Given that **Blood Pressure** has a feature importance of 0.3387, it is a "single point of failure" for the system. In production, strict validation for blood pressure format and range must be added at the API interface to prevent prediction failures due to missing or erroneous data.
- **Missing Value Strategy:** The model was trained treating "unlabeled" data as "healthy." In deployment, if users skip non-core fields, the system should default to median values or specific markers to ensure service availability, but must indicate the impact on data integrity via a Confidence Score in the returned result.

## 10.0  Limitations
While the model performs excellently on the test set, we must honestly address the following limitations before broad application:

- Insufficient Data Scale & Diversity: Only **374 sample records**. While SVM is robust with small samples, such a small dataset struggles to capture complex population distribution characteristics, carrying a risk of overfitting specific groups (e.g., specific age groups or occupations).
- **The "Stress Paradox" & Subjective Bias:** Analysis shows user self-reported "Stress Level" contributes almost zero to prediction (Feature Importance 0). This implies subjective self-reporting may have severe bias or high multicollinearity with objective physiological markers (Heart Rate, Sleep Quality), making it unreliable as an independent diagnostic basis.
- **Specific Class Confusion:** The model shows minor confusion (approx. 3 cases) distinguishing "Insomnia" from "Sleep Apnea". This is primarily because both share highly overlapping symptoms across the limited 11 features (e.g., high blood pressure, poor sleep quality), lacking features with higher discriminative power (e.g., SpO2, snoring frequency).

## 11.0  Future Work

Suggested future works to be considered are:

- Integration with Wearable Technology (IoT)**:** Transitioning from static data to continuous streams from wearables (e.g., Apple Watch, Fitbit) would enable real-time monitoring of heart rate variability, sleep stages, and activity patterns. This shift from one-time risk assessment to ongoing health tracking would validate model performance in real-world settings through ethically approved pilot studies.

- Time-Series Analysis with Deep Learning: Once continuous data streams are available, the modeling approach can evolve from static classification to Time-Series Forecasting.  Unlike our current ensemble of SVM-based approach, which analyse snapshots in time, Recurrent Neural Networks (RNNs) or Long Short-Term Memory (LSTM) networks could "remember" temporal patterns in sleep data. This would allow the system to predict sleep disorder onset before it becomes chronic, capturing deterioration trends over weeks or months.

- From Prediction to Intervention (AI Health Coach): Currently, the model outputs a predicted likelihood of disorder. The next value-added step is Prescriptive Analytics. Future versions should build a recommendation engine on top of the risk prediction. For example, if the model detects high "Work Pressure" and low "Physical Activity" as the primary contributors to a user's risk, it could automatically generate personalized, actionable advice (e.g., "Your risk of Insomnia has risen by 15%; consider a 20-minute walk to lower cortisol levels").

## References

Alshammari, T. S. (2024). "Applying Machine Learning Algorithms for the Classification of Sleep Disorders." IEEE Access **12**: 36110–36121.

LKYSPP NUS. Sleep Deprivation in Singapore: A Public Health Crisis. https://lkyspp.nus.edu.sg/gia/article/sleep-deprivation-in-singapore-a-public-health-crisis

https://www.kaggle.com/code/viniciusdiasx/predicting-sleep-disorders-with-machine-learning

https://github.com/vinnie071015/sleeping-disorder-mlops/blob/main/frontend/ui.py

https://github.com/vinnie071015/sleeping-disorder-mlops/blob/main/src/data_processor.py

## Appendix A - Data Dictionary

| Feature | Type | Description | Clinical Relevance |
|---|---|---|---|
| Person ID | Identifier | Unique identifier | - |
| Gender | Categorical (Male/Female) | Biological sex | Gender differences in sleep disorders |
| Age | Numerical (27-59) | Years | Risk increases with age |
| Occupation | Categorical (7 types) | Professional role | Occupational stress impact |
| Sleep Duration | Numerical (5.8-8.5) | Hours per night | Primary sleep metric |
| Quality of Sleep | Ordinal (1-10) | Self-reported quality | Subjective sleep assessment |
| Physical Activity Level | Ordinal (1-10) | Activity intensity | Exercise-sleep relationship |
| Stress Level | Ordinal (1-10) | Perceived stress | Stress-induced insomnia |
| BMI Category | Categorical | Under/Normal/Over/Obese | Obesity links to sleep apnea |
| Blood Pressure | Categorical | SBP/DBP ranges | Cardiovascular-sleep connection |
| Heart Rate | Numerical (65-86) | BPM | Autonomic nervous system indicator |
| Daily Steps | Numerical (3000-10000) | Steps per day | Activity metric |
| Sleep Disorder | Target (3 classes) | None/Insomnia/Sleep Apnea | Clinical diagnosis |