# Measuring Anxiety Levels with Head Motion Patterns in Severe Depression Population

Fouad Boutaleb[1,4], Emery Pierson[2], Nicolas Doudeau[4], Clémence Nineuil[4], Ali Amad[4], and Mohamed Daoudi[1,3]

[1] Univ. Lille, CNRS, Centrale Lille, Institut Mines-Télécom, UMR 9189 CRIStAL, F-59000 Lille, France
[2] LIX, École Polytechnique, IP Paris
[3] IMT Nord Europe, Institut Mines-Télécom, Univ. Lille, Centre for Digital Systems, F-59000 Lille, France
[4] Univ. Lille, Inserm, CHU Lille, U1172 - LilNCog - Lille Neuroscience & Cognition, F-59000 Lille, France

*Abstract*— Depression and anxiety are prevalent mental health disorders that frequently cooccur, with anxiety significantly influencing both the manifestation and treatment of depression. An accurate assessment of anxiety levels in individuals with depression is crucial to develop effective and personalized treatment plans. This study proposes a new noninvasive method for quantifying anxiety severity by analyzing head movements -specifically speed, acceleration, and angular displacement - during video-recorded interviews with patients suffering from severe depression. Using data from a new CALYPSO Depression Dataset, we extracted head motion characteristics and applied regression analysis to predict clinically evaluated anxiety levels. Our results demonstrate a high level of precision, achieving a mean absolute error (MAE) of 0.35 in predicting the severity of psychological anxiety based on head movement patterns. This indicates that our approach can enhance the understanding of anxiety's role in depression and assist psychiatrists in refining treatment strategies for individuals.

## I. INTRODUCTION

Depression is a prevalent mental health disorder that affects approximately 280 million people worldwide [33]. According to the Diagnostic and Statistical Manual of Mental Disorders (DSM-5), the symptoms of Clinical Depression can manifest in many various ways, including severe depression with suicidal thoughts [4]. Objective means of detecting severe depression are crucial to enable early diagnosis and better medical treatments. Recent research has shown that multiple, non-verbal, behavioral indicators of depression can be used for this objective [5], [2], [10].

However, depression remains a complex condition that varies significantly between individuals [3]. Severely depressed patients often experience additional symptoms, such as psychomotor retardation or anxiety, which can present as physical or psychological distress. Psychiatrists frequently observe these symptoms. Evaluation of these symptoms is critical in guiding treatment decisions and supporting patient recovery [8]. In other words, understanding the severity of each symptom allows clinicians to tailor treatments more effectively to each patient.

The use of behavioral indicators to provide a detailed analysis of a patient's status has not yet been fully explored. Patient examinations are heavily based on subjective

evaluations, depending on the clinician's experience and the patient's ability to communicate their symptoms [14]. This variability affects the reliability of diagnoses and may miss subtle but important differences in symptoms such as anxiety. Automatic detection algorithms based on non-verbal cues could help overcome these challenges by offering more consistent and objective evaluations, providing critical insights into managing severely depressed patients.

Wearable physiological devices, such as heart rate monitors or electrodermal activity sensors, have been used to detect several anxiety disorders.While these physiological measures provide valuable information, they are often costly and not easily accessible for widespread clinical use. Additionally, the intrusive nature of wearable devices can impact patient comfort and compliance, limiting their practicality in routine assessments.

At the same time, it has become a general fact that emotion and patients' feelings can be understood from visual media [30]. Moreover, extracting behavioral markers, such as head movements, from videos offers the advantage of being discreet and easily integrated into standard interactions, providing continuous and real-time analysis without specialized equipment.

To explore the possibility of obtaining personalized, adaptive treatments based on objective, non-intrusive markers, we captured the CALYPSO Depression Dataset. This dataset includes data from patients diagnosed with severe depression, along with detailed psychiatrist-assessed anxiety levels. Moreover, it contains informal interviews designed to reproduce daily life scenarios, such as meetings with a general practitioner.

Won et al. [32] have shown that head movements are a useful indicator of anxiety. Anxiety, whether psychological or physiologic, is a common symptom of severe depression. In particular, it intensifies the effects of depression and complicates treatment [8]. The presence of psychological anxiety is always reported in Hamilton interviews [17], [13], making this data available for our approach. We introduce a novel method to assess whether or not, measuring anxiety levels of severely depressed patients can be done using head motion information. This method is by nature, non-invasive and analyzes head movements by separating motion and non-motion sequences during the interview. Our pipeline pro-

vides objective, measurable, and interpretable data that can improve diagnostic accuracy and support more personalized treatment strategies.

Our approach emphasizes objectivity by focusing on non-depression-specific behavioral markers—head motion dynamics—that are observable even during a patient's first clinical interaction. Unlike methods requiring longitudinal tracking or disorder-specific symptom coding, we analyze motion patterns inherently linked to anxiety, ensuring applicability in real-world scenarios where clinicians may lack prior patient history. By avoiding depression-specific features, our method enhances generalizability, providing a practical tool for rapid anxiety assessment across diverse populations. It could also be applicable to other patient groups and non-clinical populations, further extending its utility. This aligns with clinical workflows, where initial interviews often serve as the primary basis for early diagnosis and treatment planning.

## II. RELATED WORK

### A. Automatic Depression Detection Using Nonverbal Cues

Clinicians have observed for a long time that severe depression is correlated with reduced physiological activities, such as monotonic tone, reduced intensity of facial expressions, or low quantity of body motion [4].

Several studies have indeed shown that modern computer science tools can be used to extract those non-verbal cues and provide an objective way to assess the level of depression. Features extracted from body gestures [20], speech patterns [9], facial expressions [11], and head movement [21] are used to automatically and accurately predict depression severity. However, the interpretability of most approaches remains unclear. Many recent approaches rely on deep neural networks [28], [29], with improved accuracy but providing limited insights for clinicians.

Most studies on detecting depression focus on separating severely depressed and healthy populations but lack interpretability. Gahalawat et al. [15] addressed this issue by proposing an interpretable approach using head motion patterns for binary classification. However, while their method enhances explainability, it remains focused on distinguishing between groups rather than evaluating the precise state of depressed patients.

### B. Automatic Anxiety Detection

Similar to depression, research has shown that several anxiety disorders can be automatically detected using non-verbal cues [24]. Most studies have primarily focused on physiological signals—such as heart rate variability, skin conductance, and cortisol levels—using wearable devices [23], [19].

In [27], the authors extract facial features, such as face and mouth motion, to predict whether patients are anxious or relaxed. This analysis was later extended in [16]. Mo et al. [26] propose using facial cues to accurately detect anxiety and distress in a non-intrusive manner. However, their feature extraction process relies on deep learning and is therefore non-interpretable.

In our approach, anxiety is considered a symptom of depression, whereas most of the cited works treat anxiety as an independent illness. It remains unclear whether these methods would be effective in detecting anxiety in depressed patients.

### C. Tracking behavior with head motion

Head movements are significant nonverbal indicators for behavioral analysis and are well studied, as they are easy to track in virtual reality environments [25] or dyadic interactions [31]. In virtual reality settings, head movements serve as valuable features for analyzing social interactions [18], emotional states [34], and simulation sickness in virtual environments [6].

By extracting head motion dynamics from videos of structured Hamilton interviews, Kacem et al. [21] classified depression severity, finding that depressed individuals exhibit less head movement. Dibeklioglu et al. [12] combined head movements with facial dynamics and vocal prosody for depression detection, noting differences in nodding frequency and amplitude between depressed and non-depressed populations. Finally, head movements have also been shown to be valuable for anxiety prediction [32].

We summarize the main contributions of our work below:

1) We introduce the CALYPSO dataset, a longitudinal study of clinical depression. In particular, we propose to use the videos of informal interviews of the study to analyze anxiety in severe depression using head movements extracted from the videos.
2) We propose segmenting videos into head-moving and non-moving phases. This approach allows us to extract a more comprehensive set of head motion features for our analysis. We further select the most significant features and train a regression model to predict psychological anxiety.
3) We validate our method on the CALYPSO dataset and demonstrate its effectiveness in daily life settings by applying it to informal interviews with depressed patients. For the first time, we establish a link between anxiety in severe depression and objectively measurable nonverbal cues, based on head motion characteristics extracted from pose angular displacements.

## III. METHODOLOGY

In this section, we outline our methodology. Our approach is based on videos of informal interviews from the CALYPSO dataset. We divide our approach into several steps. First, we automatically extracted the head pose and detailed its motion from the videos. Then, we extracted statistical features from the poses and motion sequences. Our final step is a trainable pipeline to select features and train a linear, interpretable model to predict physical anxiety from the selected features.
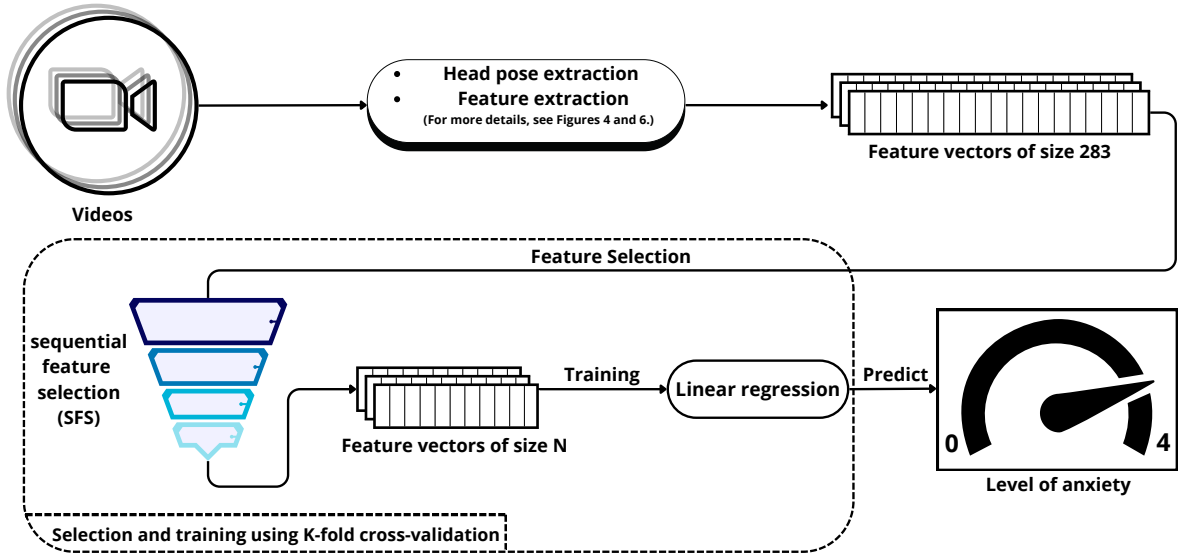
Fig. 1. Overview of our proposed pipeline. We apply the same process to all videos of the informal interviews of the CALYPSO dataset. We first extract the head pose and its motion (speed and acceleration) automatically. We then apply statistical feature analysis to extract a feature vector of size 283 for each video. Finally, we apply a cross-validated approach to select features and train a linear model to accurately regress psychological anxiety levels.
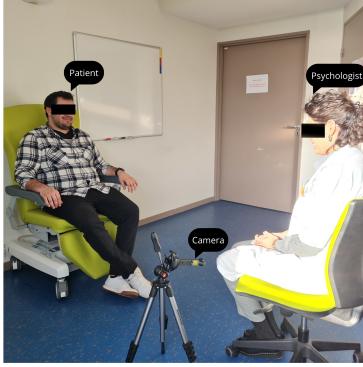


Fig. 2. Interview Room Setup for the Calypso Depression Dataset.

### A. Head pose and motion extraction

**Pose Extraction**

To track head orientation in 3D space throughout the video, we used MediaPipe software [22]. We captured the head orientation using Euler angles—pitch, yaw, and roll—representing the different axes of rotation.

As shown in the figure 3, the three primary axes of rotation are:

- **Pitch** ($\theta$): Rotation around the x-axis (nodding up and down).
- **Yaw** ($\phi$): Rotation around the y-axis (turning left and right).
- **Roll** ($\psi$): Rotation around the z-axis (tilting the head side to side).

This process resulted in a time series of angles that describe the head's orientation throughout the interview, as shown in Figure 4.
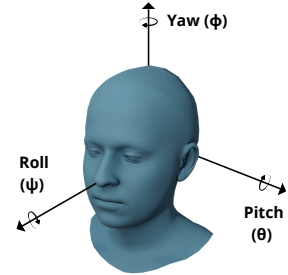
**Angular Velocity Calculation**



Fig. 3. Diagram of head motion axes—pitch, roll, and yaw—used in our analysis.

To measure head motion, we calculated the angular velocity for yaw, pitch, and roll. Let t be a time step and $t+\Delta t$ be the following step. As the pose can be described as a rotation matrix (computed from the yaw pitch and roll angles), we compute the derivative of the rotation matrix $R_t$.

This derivative is computed using the following formula:

$$\omega(t) \approx \frac{1}{\Delta t} \left( R_{t+\Delta t} R_t^T - I \right).$$

The product $R_{t+\Delta t} R_t^T$ is the relative rotation between two consecutive poses, and we measure its deviation from the identity matrix $I$. The resulting matrix $\omega(t)$ is a skew-symmetric matrix, from which we can recover the angular velocity across yaw, pitch, and roll axes:

$$\omega(t) = \begin{bmatrix} 0 & -\omega_z & \omega_y \\ \omega_z & 0 & -\omega_x \\ -\omega_y & \omega_x & 0 \end{bmatrix}$$

where $\omega_x$ is the angular velocity around the pitch-axis, $\omega_y$ is the angular velocity around the yaw-axis, and $\omega_z$ is the angular velocity around the roll-axis.
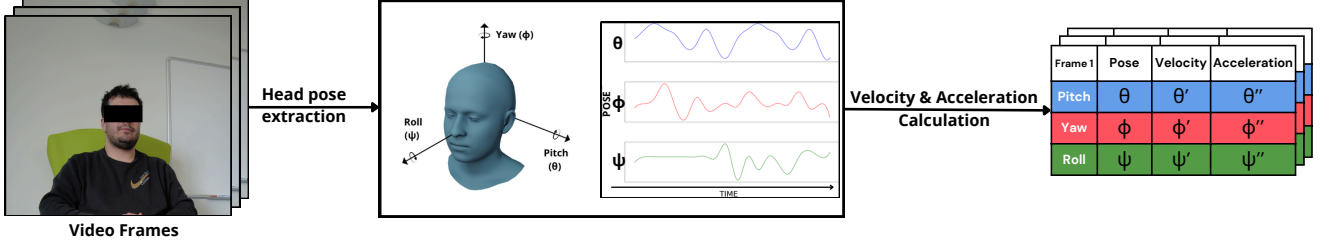
Fig. 4.    Illustration of the head pose and motion extraction.

**Acceleration Calculation**

Since the skew-symmetric matrix space is linear, calculating the angular acceleration is simply derived from:

$$\dot{\omega}(t) \approx \frac{\omega_{t+1} - \omega_t}{\Delta t}$$

where $\omega_t$ and $\omega_{t+1}$ are the angular velocities at consecutive timestamps $t$ and $t + 1$, and $\Delta t$ is the time interval between the measurements.

The pitch, yaw, and roll acceleration are defined as $\dot{\omega}_x, \dot{\omega}_y, \dot{\omega}_z$.

*B. Motion segmentation*

We observed that statistical features extracted from the full interview sequences are not discriminative enough to provide reliable predictions (see Table 1). Moreover, we noted that the head pose of the patient alternates between two states: **moving** (the patient is changing position on the chair) and **steady** (the patient is on a stable position and has limited motion). This motivated us to segment interviews in moving and steady sequences.

To provide a flexible yet straightforward method for achieving this, the velocity data was clustered into two groups using a Gaussian Mixture Model (GMM). The two clusters are effective at classifying head movement between the intuitive "moving" and "steady" states.

This approach can be generalized across different patients, making the classification robust for various head movement behaviors. The whole process is illustrated in Figure 5.
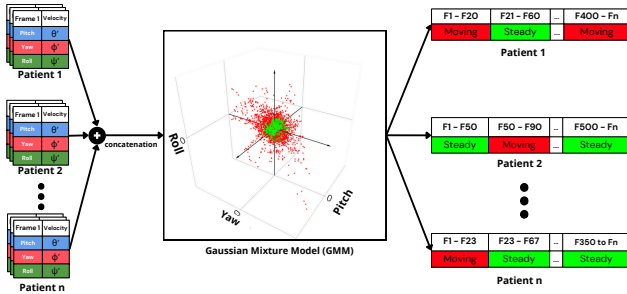


Fig. 5.    We apply a Gaussian Mixture Model (GMM) to cluster head rotational velocities (pitch, yaw, roll) into "moving" and "steady" states, segmenting interviews into sequences (e.g., F21-F60, representing frames 21 to 60). The plot illustrates the velocity profiles for each axis.

*C. Feature Extraction*

Using the resulting clustering from the Gaussian Mixture Model (GMM), the interview videos were segmented into moving and steady sequences. For each segment, head movements were analyzed in terms of pitch, yaw, and roll, along with their velocities and accelerations. This resulted in a total of nine core features: three rotational angles (pitch, yaw, roll), their velocities, and their accelerations. (9 features across 2 clusters)

To comprehensively capture the characteristics of head movement, we extracted three types of statistical features:

- **Global statistical features:** For a general overview, we grouped the "moving" and "stable" segments into two subsets and calculated summary statistics, including Mean, Median, Range, Median Absolute Deviation (MAD), Skewness, Kurtosis, and Standard Deviation for each feature. This resulted in a total of 7 statistics for 9 features across 2 clusters (moving and stable). ($7 \times 9 \times 2$)
- **Sequence-Level Features:** We then separated each moving and steady sequence and analyzed them individually. We then computed the following 7 statistics: Mean, Median, Range, MAD, Skewness, Kurtosis, and Standard Deviation of pitch, yaw, roll, velocities, and accelerations. These features were then averaged across all sequences in each group (moving or steady) ($7 \times 9 \times 2$). Additionally, we included the sum of absolute values for all speeds and accelerations (cumulative displacement, cumulative acceleration). ($1 \times 6 \times 2$).
- **Temporal Features:** For the temporal analysis, we focused on the duration of each movement or stable segment. We calculated the mean, median, standard deviation, skewness, range of durations, and the ratio of time spent in each state (moving vs. stable). Additionally, we computed the number of transitions per minute, which reflects how often the subject switched between movement and stillness. ($6 \times 2 + 1$)

This process resulted in a feature vector of size 283 ( 126 global, 144 sequence-level, and 13 temporal) for each video.

*D. Feature Selection and Model Development*

To reduce the risk of overfitting given the high number of features, we applied a selection method that reduces the number of features. Such approaches have shown useful
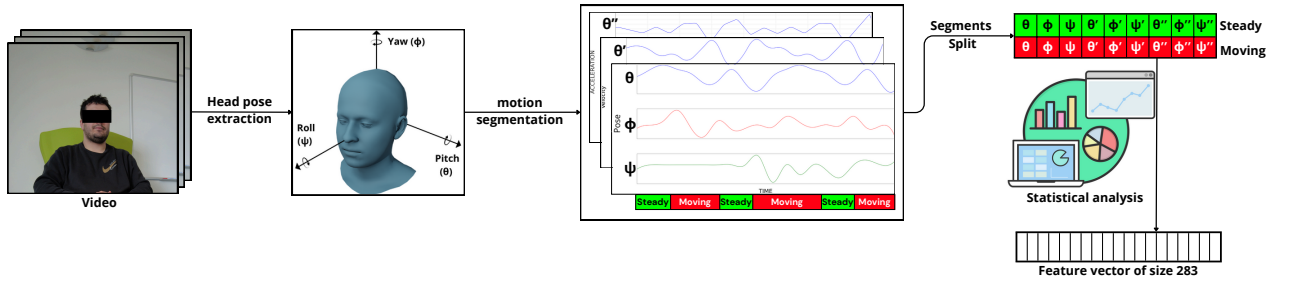
Fig. 6. Illustration of the full feature extraction process.

in extracting interpretable features for depression assessment [1], [7].

*Correlation Filtering* To ensure that the model was not impacted by multicollinearity, and to simplify the process, we first computed the correlation between all extracted features. Any features with a correlation coefficient greater than ±0.8 were removed (with a preference for non-derivative features). This step reduced redundancy and eliminated high correlation in features that could affect the reliability of the model. After applying this filtering process, we retained a refined set of **96 features**.

We then employed Sequential Feature Selection (SFS), which allowed us to identify the most relevant features for our task. Sequential Feature Selection is a feature selection technique that iteratively adds or removes features based on their contribution to the model's performance. In our approach, we explored two variations of SFS:

*Sequential Selection by Exclusion:* This method sequentially removes the least significant feature based on a predefined performance criterion (e.g., Mean Squared Error (MSE)). At each iteration, the feature whose exclusion results in the least degradation or the most improvement in model performance is removed. This process continues until no features are left. We then select the set of features that provided the best performance during the process.

*Sequential Backward Floating Selection (Inclusion/Exclusion):* This approach adds a conditional inclusion step to backward selection. After removing a feature, the algorithm evaluates whether reintroducing any of the previously excluded features can enhance the model's performance. This mechanism thus dynamically adjusts the feature subset, which allows recovering from suboptimal exclusions and identifying a more optimal set of features.

### E. Regression model

We used a linear regression model to allow for interpretation of the model behavior. We used the Lasso regularization as it consistently yielded the best results. All steps using training data as input were included in the cross-validation process to ensure our results were not overfitted or exhibiting spurious correlations.

### IV. RESULTS

In this section, we present the outcomes of our regression analysis aimed at predicting anxiety levels based on head
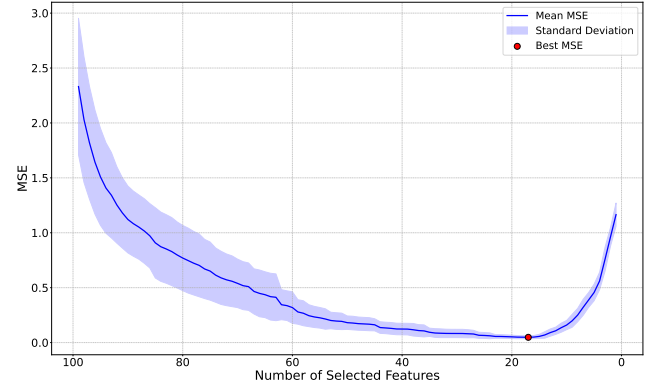


Fig. 7. The plot shows the cross-validation (CV) scores during feature selection across 10 folds. The blue line represents the mean CV score as the number of selected features increases, with the shaded area indicating the standard deviation. The x-axis shows the number of selected features, and the y-axis shows the mean squared error (MSE). The red marker indicates the mean of the lowest CV scores across all folds.

movement features. We evaluate multiple machine learning models under different feature selection processes and assess the impact of incorporating motion segmentation. We finally provide an interpretation of our model behavior and features that are shown to be linked with the presence or absence of anxiety in severe depression.

### A. CALYPSO dataset

*Patient Selection*
Patients admitted to the hospital undergo standard diagnostic evaluations with an attending psychiatrist. A clinician then conducts a first examination to determine if the patient's history and symptoms align with the DSM-5 (Diagnostic and Statistical Manual of Mental Disorders) conditions for severe clinical depression [4]. Patients who meet the inclusion criteria are proposed for a more in-depth interview, contributing directly to the CALYPSO depression dataset. In total, **32 patients** meeting these criteria were included in the study. The CALYPSO clinical trial has been reviewed and approved by the Ethics Committee under approval number 2022-A01160-43, ensuring adherence to ethical standards and guidelines. All patients provided written informed consent before participation in the study.

The study participants were predominantly French nationals (ethnicity data were not collected), with an equal gender

distribution (50% male, 50% female).

*Interview Process*

Selected patients participated in structured clinical interviews conducted in a controlled environment, as shown in Figure 2. The interview was divided into two distinct phases: an initial informal segment that consists of a casual conversation between the psychologist and the patient, lasting only a few minutes, during which the patient is asked non-medical questions.

After the informal conversation, the clinician conducts a structured interview aimed at evaluating the Hamilton Depression Rating Scale (HDRS) for the patient. This assessment specifically includes measuring psychological anxiety. After the interview, the clinician assigns a psychological anxiety score ranging from 0 (no anxiety) to 4 (high anxiety), based on the patient's responses and observed behaviors, as shown in Figure 8.
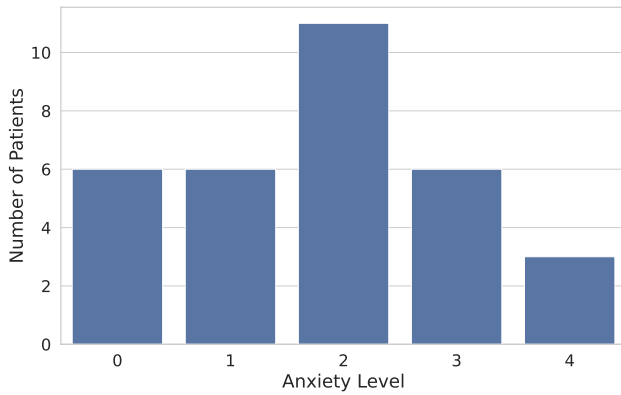


Fig. 8. Distribution of Anxiety Levels in the CALYPSO Dataset.

To assess if the head motion patterns can be used in daily life interviews, such as discussions with a general practitioner, we conducted our analysis on the video recordings from the informal part of the interview.

*B. Experimental setup*

To identify the most effective set of features and minimize the risk of overfitting, we employed Sequential Feature Selection (SFS) in combination with 10-fold cross-validation. We evaluated multiple machine learning models and fine-tuned the alpha parameter to determine the optimal configuration for our dataset. In each cross-validation fold, SFS was applied to the training data (comprising 9 folds) to select the best-performing features, resulting in 10 distinct feature lists. We then consolidated these results by selecting features that appeared in at least five of the ten lists, ensuring that only the most consistently important features were retained. This approach enhanced the model's robustness and generalizability by focusing on reliable predictors. Finally, we trained the final model using this refined set of features, which streamlined the model and improved its performance on unseen data.

*C. Evaluation Metrics*

The regression performance was evaluated using the mean absolute error (MAE);

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i| \qquad (1)$$

and the coefficient of determination (R2 score),:

$$R2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2} \qquad (2)$$

where, $n$ is the number of observations, $y_i$ is the actual value, $\hat{y}_i$ is the predicted value, and $\bar{y}$ is the mean of the actual values.

To further assess the practical applicability of our regression model, we converted the regressed predictions into discrete values. Each predicted value is converted to the closest integer. We then report the classification accuracy based on the predicted anxiety level.

*D. Psychological Anxiety Level Prediction*

For predicting psychological anxiety levels, the best-performing model was Lasso regression, achieving a Mean Absolute Error (MAE) of 0.31 and a coefficient of determination ($R2$) of 0.87 with only 14 features needed (see Figure 9).
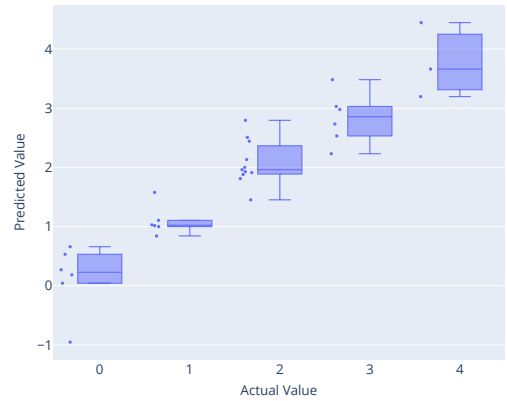


Fig. 9. This figure illustrates the actual versus predicted psychological anxiety levels for the best Lasso model, each point represents a patient

*E. Impact of segmentation*

To evaluate the significance of the interview motion segmentation in our feature extraction process, we compared it to a baseline model without the segmentation. The baseline approach involved extracting features directly from the raw head motion data without segmentation. We computed the same statistical measures (mean, median, range, skewness, kurtosis, standard deviation) for the head angles (pitch, yaw, roll), velocities, and accelerations, for a total of 54 features.

As shown in Table I, our approach significantly improves the predictive performance of the model compared to the baseline.

TABLE I

REGRESSION RESULTS FOR PSYCHOLOGICAL ANXIETY LEVEL

| Selection Process | Features | MAE | R2 | Accuracy |
|---|---|---|---|---|
| Full model (ours) | 14 | 0.31 | 0.87 | 0.75 |
| Without GMM | 16 | 0.90 | 0.10 | 0.44 |

*F. Comparison of Feature Selection Methods*

We also compared different feature selection methods to evaluate their impact on model performance. Specifically, we examined exclusion-only selection, inclusion and exclusion (I/E) selection, and no selection at all. Table II summarizes the results.

TABLE II

COMPARISON OF FEATURE SELECTION METHODS FOR PSYCHOLOGICAL ANXIETY LEVEL

| Selection Process | No. of Features | MAE | R2 | Accuracy |
|---|---|---|---|---|
| Exclusion Only | 12 | 0.48 | 0.76 | 0.50 |
| Inclusion and Exclusion (I/E) | 14 | 0.31 | 0.87 | 0.75 |
| I/E - Without GMM | 16 | 0.90 | 0.10 | 0.44 |
| No selection | 283 | 0.98 | -0.05 | 0.31 |

The results indicate that the inclusion and exclusion (I/E) feature selection method yields the best performance. Notably, the exclusion only process results in fewer features than the inclusion and exclusion strategy. Moreover, the Lasso $L_1$ penalty alone is not sufficient to select the features without overfitting or satisfying performance (R2 score is almost zero, meaning that the model always predicts the mean value).

*G. Classification Performance Analysis*

To provide a qualitative analysis of the classification, we present the confusion matrix of our model in Figure 13.



Fig. 10. Confusion Matrix: Classification results by grouping continuous anxiety level predictions into classes with a tolerance of ±0.5 units.

We observe that all wrong predictions differ by one, which we find acceptable given the inherent variability in psychiatrist ratings. This result suggests that our model's performance is comparable to human-level error margins in clinical settings.

*H. Interpretation of results*

Figure 11 presents the best model coefficients derived from the regression analysis predicting anxiety levels based on head movement data. Each bar represents the importance or weight of a specific feature, with the feature names listed on the y-axis and their respective coefficient values on the x-axis. The key observations include:

- **Global — Pitch Degree — Median — Steady**: This feature shows the highest negative coefficient, indicating that steady pitch movements (up and down head movements) are inversely linked to anxiety levels. In other words, when the median pitch degree is low (indicating the head is raised), anxiety levels tend to be higher.

- **Temporal — Skewness — Moving**: The second-highest coefficient (positive) indicates that patients with high anxiety tend to spend longer periods in sequences where they are moving. Skewness highlights the imbalance in time spent during movement. In other words, patients with a high level of physical anxiety have an irregular duration of motion, and alternate between long moments of motion, and shorter ones, compared to more stability for non anxious patients.

- **Temporal — Visits per Minute**: This feature measures the number of transitions between moving and steady periods per minute during the observation period. It correlates negatively with anxiety levels, which might seem counterintuitive at first. However, this suggests that more anxious patients tend to either stay moving for long periods or remain still for extended durations, resulting in fewer transitions between motion and steady states. By combining information from other features, we deduce that higher anxiety levels are associated with patients staying in motion for longer stretches, while less anxious patients frequently alternate between moving and stopping.

The density plots in Figure 12 illustrate how these key features vary across different anxiety levels, offering a more detailed insight into the trends that correspond with the model's predictions.

Overall, the model suggests that stable head motion is linked to lower anxiety, while unstable head motion, particularly with longer and more irregular motion sequences, is associated with higher anxiety.

## V. LIMITATIONS AND FUTURE WORK

*A. Limitations*

While our method demonstrates strong performance in predicting psychological anxiety levels using head motion patterns, it is less effective in assessing somatic anxiety. We applied the same technique to provide a regression model for somatic anxiety levels. However, the regression models yielded lower predictive accuracy, with a Mean Absolute Error (MAE) of **0.47** and an R² score of **0.53**. With these results, the large error obtained does not allow for a clear differentiation between the affected and non-affected populations.
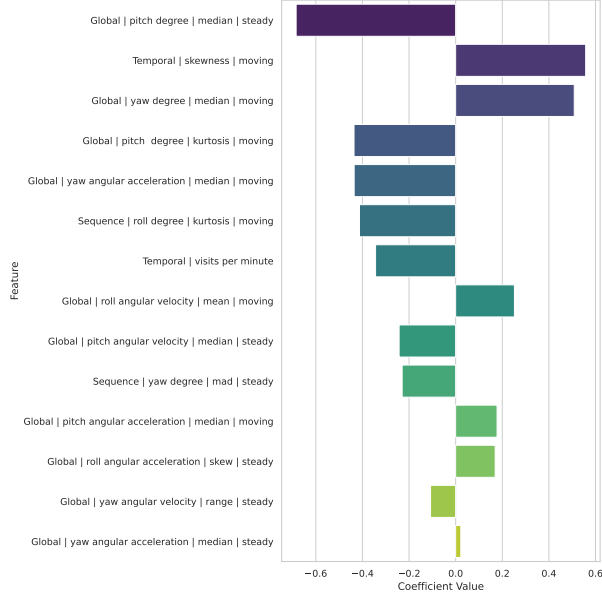
Fig. 11. This plot illustrates the coefficients of the Lasso regression model for each feature in the best model.
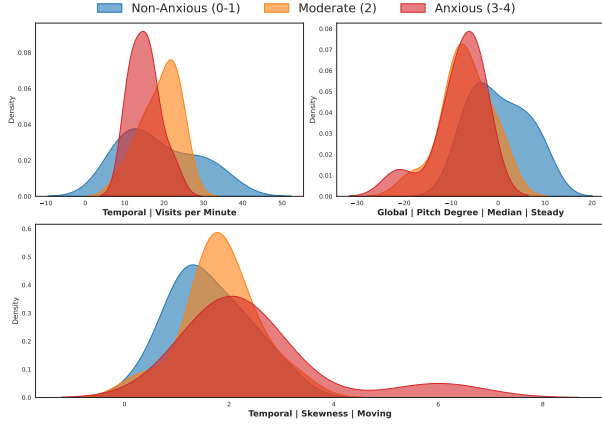


Fig. 12. Density distributions of key features across different anxiety score groups. The blue, orange, and red areas represent the density estimates for individuals with anxiety scores of 0-1 (non-anxious), 2 (moderately anxious), and 3-4 (highly anxious), respectively.

Somatic anxiety manifests through physical symptoms such as muscle tension, restlessness, and other bodily sensations that may not be fully captured by analyzing head movements alone. We believe those manifestations often involve subtle physiological changes or whole-body movements that require additional modalities to detect accurately. Incorporating other physiological or behavioral cues may be necessary to comprehensively assess this subtype of anxiety in individuals with severe depression.

*B. Future Works*

Several avenues of research are possible to validate and expand upon this work.

**Multimodal Analysis:** This study focuses on head motion patterns. However, CALYPSO dataset contains full videos of patients. Integrating other features, from facial expressions, body gestures or speech patterns could help provide measurements of depression symptoms. Combining multimodal data may improve the predictive accuracy of the model and offer deeper insights into the behavioral manifestations of anxiety in depression.

**Longitudinal Studies:** In this study, our goal was to develop a model that accurately predicts anxiety levels in depressed patients based on their first clinical interview. However, future work will explore applying this approach to longitudinal data to monitor anxiety evolution throughout treatment. Investigating whether changes in head motion patterns correlate with treatment response could provide valuable insights for adjusting therapeutic strategies over time.

## VI. CONCLUSION

In this study, we introduced a novel, non-invasive method for quantifying anxiety severity in patients with severe depression by analyzing head motion patterns during clinical interviews. We propose also a new depression dataset named CALYPSO, introduced in this paper, which contains video data of depressed patients, from which we extracted features related to head motion. Moreover, our new approach, separating moving and non-moving segments of the interview, allows us to extract more valuable features for the analysis. We demonstrated that we can train an interpretable model based on the selected features for predicting the anxiety level of depressed patients, achieving a Mean Absolute Error (MAE) of 0.31 and an $R^2$ of 0.87.

These results suggest that head motion patterns can serve as reliable, objective indicators of anxiety severity in individuals with severe depression. By providing an automated and quantifiable assessment tool, our approach has the potential to assist psychiatrists in making more informed decisions regarding diagnosis and treatment planning. This method enhances the understanding of anxiety's role in depression and contributes to more personalized and effective interventions for patients suffering from both conditions.

## VII. APPENDIX

*A. Confusion matrix of somatic Anxiety*

We show in Figure 13 the confusion matrix of our approach regression model for somatic anxiety. Notably, the moderately anxious class is hardly well predicted compared to the results of psychological anxiety in the main paper, supporting the need for more non-verbal features to assess the somatic anxiety symptom.

*B. Supplementary results for prediction of psychological anxiety*

We provide in Table III, a detailed ablation of each pipeline parameter. In particular, we show that the Lasso model is the best regularization for regressing psychological anxiety.

TABLE III

FULL ABLATION STUDY OF OUR MODEL FOR PSYCHOLOGICAL ANXIETY. THE BEST MODEL IS THE LASSO MODEL WITH I/E SELECTION PROCESS.

| Selection Process for SFS | Model | Alpha[a] | Number of Features Selected[b] | MAE[c] | R2[c] |
|---|---|---|---|---|---|
| Exclusion Only | Ridge | 1.0 | 23 | **0.35** | **0.87** |
| | | 0.1 | 16 | 0.49 | 0.76 |
| | Lasso | 0.01 | 12 | 0.48 | 0.76 |
| | Linear Regression | N/A | 23 | 0.53 | 0.64 |
| | ElasticNet | 0.01 | 13 | 0.59 | 0.65 |
| Inclusion and Exclusion | Ridge | 1.0 | 18 | 0.46 | 0.78 |
| | | 0.1 | 16 | 0.49 | 0.76 |
| | Lasso | 0.1 | 6 | 0.63 | 0.61 |
| | | 0.01 | 14 | **0.31** | **0.87** |
| | | 0.001 | 16 | 0.63 | 0.60 |
| | Linear Regression | N/A | 13 | 0.72 | 0.45 |
| | ElasticNet | 0.01 | 9 | 0.57 | 0.68 |
| | | 0.001 | 14 | 0.60 | 0.67 |
| Inclusion and Exclusion (Without GMM) | Ridge | 1.0 | 17 | 0.80 | 0.35 |
| | | 0.1 | 24 | **0.67** | **0.45** |
| | Lasso | 0.01 | 16 | 0.90 | 0.10 |
| | Linear Regression | N/A | 30 | 1.6 | -2.45 |
| | ElasticNet | 0.01 | 7 | 0.76 | 0.42 |



Fig. 13. Confusion Matrix obtained using the best model from Table IV. Classification results are shown by grouping continuous anxiety level predictions into classes with a tolerance of ±0.5 units.

## C. Supplementary results for prediction of somatic anxiety

We provide in Table IV, a similar ablation for somatic anxiety. Notably, the Lasso model fails to provide similar accuracy. Moreover, no model is able to provide similar results as for psychological anxiety.

REFERENCES

[1] S. Alghowinem, T. Gedeon, R. Goecke, J. F. Cohn, and G. Parker. Interpretation of depression detection models via feature selection methods. *IEEE Transactions on Affective Computing*, 14(1):133–152, 2023.

[2] S. Alghowinem, R. Goecke, M. Wagner, G. Parker, and M. Breakspear. Eye movement analysis for depression detection. In *2013 IEEE International Conference on Image Processing*, pages 4220–4224. IEEE, 2013.

[3] K. Allsopp, J. Read, R. Corcoran, and P. Kinderman. Heterogeneity in psychiatric diagnostic classification. *Psychiatry research*, 279:15–22, 2019.

[4] D. American Psychiatric Association, D. American Psychiatric Association, et al. *Diagnostic and statistical manual of mental disorders: DSM-5*, volume 5. American psychiatric association Washington, DC, 2013.

[5] N. V. Babu and E. G. M. Kanaga. Sentiment analysis in social media data for depression detection using artificial intelligence: a review. *SN computer science*, 3(1):74, 2022.

[6] J. N. Bailenson and N. Yee. A longitudinal study of task performance, head movements, subjective report, simulator sickness, and transformed social interaction in collaborative virtual environments. *Presence: Teleoperators and Virtual Environments*, 15(6):699–716, 2006.

[7] M. Bilalpur, S. Hinduja, L. A. Cariola, L. B. Sheeber, N. Alien, L. A. Jeni, L.-P. Morency, and J. F. Cohn. Multimodal feature selection for detecting mothers' depression in dyadic interactions with their adolescent offspring. In *2023 IEEE 17th International Conference on Automatic Face and Gesture Recognition (FG)*, pages 1–8, 2023.

[8] K. W. Choi, Y.-K. Kim, and H. J. Jeon. Comorbid anxiety and depression: clinical and conceptual consideration and transdiagnostic treatment. *Anxiety disorders: Rethinking and understanding recent discoveries*, pages 219–235, 2020.

[9] N. Cummins, S. Scherer, J. Krajewski, S. Schnieder, J. Epps, and T. F. Quatieri. A review of depression and suicide risk assessment using speech analysis. *Speech communication*, 71:10–49, 2015.

[10] M. Deshpande and V. Rao. Depression detection using emotion artificial intelligence. In *2017 international conference on intelligent sustainable systems (iciss)*, pages 858–862. IEEE, 2017.

[11] H. Dibeklioğlu, Z. Hammal, and J. F. Cohn. Dynamic multimodal measurement of depression severity using deep autoencoding. *IEEE journal of biomedical and health informatics*, 22(2):525–536, 2017.

[12] H. Dibeklioglu, Z. Hammal, Y. Yang, and J. F. Cohn. Multimodal detection of depression in clinical interviews. In *ICMI*, pages 307–310. ACM, 2015.

[13] M. B. First and M. Gibbon. The structured clinical interview for dsm-iv axis i disorders (scid-i) and the structured clinical interview for dsm-iv axis ii disorders (scid-ii). 2004.

[14] C. FitzGerald and S. Hurst. Implicit bias in healthcare professionals: a systematic review. *BMC medical ethics*, 18:1–18, 2017.

[15] M. Gahalawat, R. Fernandez Rojas, T. Guha, R. Subramanian, and R. Goecke. Explainable depression detection via head motion patterns. In *Proceedings of the 25th International Conference on Multimodal Interaction*, ICMI '23, page 261–270, New York, NY, USA, 2023. Association for Computing Machinery.

[16] G. Giannakakis, M. Pediaditis, D. Manousos, E. Kazantzaki, F. Chiarugi, P. G. Simos, K. Marias, and M. Tsiknakis. Stress and anxiety detection using facial cues from videos. *Biomedical Signal Processing and Control*, 31:89–101, 2017.

[17] M. Hamilton. A rating scale for depression. *Journal of neurology, neurosurgery, and psychiatry*, 23(1):56, 1960.

[18] F. Herrera and J. N. Bailenson. Virtual reality perspective-taking at scale: Effect of avatar representation, choice, and head movement on prosocial behaviors. *new media & society*, 23(8):2189–2209, 2021.

[19] B. A. Hickey, T. Chalmers, P. Newton, C.-T. Lin, D. Sibbritt, C. S. McLachlan, R. Clifton-Bligh, J. Morley, and S. Lal. Smart devices and

TABLE IV

FULL ABLATION STUDY OF OUR MODEL FOR PSYCHOLOGICAL ANXIETY. THE BEST MODEL IS THE ELASTICNET MODEL WITH EXCLUSION SELECTION PROCESS

| Selection Process for SFS | Model | Alpha[a] | Number of Features Selected[b] | MAE[c] | R2[c] |
|---|---|---|---|---|---|
| Exclusion Only | Ridge | 1.0 | 16 | 0.76 | 0.38 |
| | | 0.1 | 22 | 0.46 | 0.65 |
| | Lasso | 0.01 | 9 | 0.61 | 0.21 |
| | Linear Regression | N/A | 24 | 1.16 | -1.66 |
| | ElasticNet | 0.01 | 10 | **0.40** | **0.69** |
| Inclusion and Exclusion | Ridge | 1.0 | 15 | **0.47** | **0.63** |
| | | 0.1 | 13 | 0.50 | 0.59 |
| | Lasso | 0.01 | 9 | 0.59 | 0.26 |
| | | 0.001 | 14 | 0.47 | 0.53 |
| | Linear Regression | N/A | 12 | 0.49 | 0.56 |
| | ElasticNet | 0.01 | 7 | 0.52 | 0.47 |
| Inclusion and Exclusion (Without GMM) | Ridge | 1.0 | 9 | 0.67 | 0.19 |
| | | 0.1 | 20 | 0.79 | -0.06 |
| | Lasso | 0.01 | 7 | **0.64** | **0.29** |
| | Linear Regression | N/A | 32 | 2.67 | -12.07 |
| | ElasticNet | 0.01 | 11 | 0.69 | 0.17 |

wearable technologies to detect and monitor mental health conditions and stress: A systematic review. *Sensors*, 21(10):3461, 2021.

[20] J. Joshi, R. Goecke, G. Parker, and M. Breakspear. Can body expressions contribute to automatic depression analysis? In *2013 10th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, pages 1–7. IEEE, 2013.

[21] A. Kacem, Z. Hammal, M. Daoudi, and J. F. Cohn. Detecting depression severity by interpretable representations of motion dynamics. In *FG*, pages 739–745. IEEE Computer Society, 2018.

[22] Y. Kartynnik, A. Ablavatski, I. Grishchenko, and M. Grundmann. Real-time facial surface geometry from monocular video on mobile gpus. *arXiv preprint arXiv:1907.06724*, 2019.

[23] E. Lazarou and T. Exarchos. Predicting stress levels using physiological data: Real-time stress prediction models utilizing wearable devices. *AIMS Neuroscience*, 11:76–102, 04 2024.

[24] L. K. S. d. Lima, E. D. B. d. Assis, N. Torro, et al. Facial expressions and eye tracking in individuals with social anxiety disorder: a systematic review. *Psicologia: Reflexão e Crítica*, 32:9, 2019.

[25] P. Lindner. Better, virtually: the past, present, and future of virtual reality cognitive behavior therapy. *International Journal of Cognitive Therapy*, 14(1):23–46, 2021.

[26] H. Mo, Y. Li, P. Han, X. Liao, W. Zhang, and S. Ding. Sff-da: Spatiotemporal feature fusion for nonintrusively detecting anxiety. *IEEE Transactions on Instrumentation and Measurement*, 2023.

[27] M. Pediaditis, G. Giannakakis, F. Chiarugi, D. Manousos, A. Pampouchidou, E. Christinaki, G. Iatraki, E. Kazantzaki, P. G. Simos, K. Marias, et al. Extraction of facial features as indicators of stress and anxiety. In *2015 37th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pages 3711–3714. IEEE, 2015.

[28] S. Song, L. Shen, and M. Valstar. Human behaviour-based automatic depression analysis using hand-crafted statistics and deep learned spectral features. In *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*, pages 158–165. IEEE, 2018.

[29] Y. Suhara, Y. Xu, and A. Pentland. Deepmood: Forecasting depressed mood based on self-reported histories via recurrent neural networks. In *Proceedings of the 26th International Conference on World Wide Web*, pages 715–724, 2017.

[30] J. Z. Wang, S. Zhao, C. Wu, R. B. Adams, M. G. Newman, T. Shafir, and R. Tsachor. Unlocking the emotional world of visual media: An overview of the science, research, and impact of understanding emotion. *Proceedings of the IEEE*, 111(10):1236–1286, 2023.

[31] A. S. Won, J. N. Bailenson, and J. H. Janssen. Automatic detection of nonverbal behavior predicts learning in dyadic interactions. *IEEE Transactions on Affective Computing*, 5(2):112–125, 2014.

[32] A. S. Won, B. Perone, M. Friend, and J. N. Bailenson. Identifying anxiety through tracked head movements in a virtual classroom. *Cyberpsychology, Behavior, and Social Networking*, 19(6):380–387, 2016.

[33] World Health Organization. Depression fact sheet. https://www.who.int/news-room/fact-sheets/detail/depression, 2023.

[34] T. Xue, A. E. Ali, G. Ding, and P. Cesar. Investigating the relationship between momentary emotion self-reports and head and eye movements in hmd-based 360 vr video watching. In *Extended abstracts of the 2021 CHI conference on human factors in computing systems*, pages 1–8, 2021.