# Discussing Solutions to the Data Imbalance Problem in Emotion Recognition

**Junwei Chen**

*Maynooth International Engineering College, Fuzhou University, Fuzhou, China*
*832203214@fzu.edu.cn*

*Abstract.* Emotion recognition technology has been widely used in human-computer interaction, medical health and other fields. However, in practical applications, emotion datasets often have class imbalance problems, which lead to the model being seriously biased towards the majority class, significantly reducing the recognition accuracy and reliability of minority emotion classes. This paper focuses on comparing and analyzing methods such as ESC-GAN generative data augmentation technology, DER-GCN dialogue and event relationship perception graph model, and MultiEMO multimodal fusion framework to solve the problem of imbalanced emotion recognition categories, and explores the innovations and limitations in multiple scenarios. These methods compensate for minority emotions from different angles: for example, MultiEMO significantly improves the ability to classify minority emotions through cross-modal attention mechanism and weighted contrast loss, which can not only be applied to detect the psychological emotions of patients in the medical health field, but also help to provide support for fine-grained emotion classification in security scenarios. Experimental results show that these solutions significantly improve the accuracy and F1 value of emotion recognition, especially in extremely unbalanced categories. This paper provides a systematic reference for the selection of technology for high-value scenarios such as medical monitoring and intelligent security, promotes the interdisciplinary collaborative development in the field of emotional computing, and accelerates the application transformation of this technology in practice.

*Keywords:* Emotion Recognition, Data Imbalance, Data Augmentation, Loss Optimization, Multimodal Fusion.

## 1. Introduction

Emotion recognition technology is widely used in fields such as medical monitoring and human-computer interaction, but data often poses a major challenge in the form of class imbalance [1]; this class imbalance can significantly degrade model performance, especially for minority classes that often represent the most critical situations in real-world applications. For example, in real-world scenarios, minority emotion samples such as "fear" and "disgust" are scarce, causing the model to favor the majority class (e.g., "happiness"), seriously affecting recognition accuracy. Solving this type of data imbalance problem is crucial to improving the practicality of the technology.

To address this data imbalance problem, scholars have proposed solutions mainly from the aspects of sampling strategy, data enhancement, and loss function optimization. In the research based on a sampling strategy, Hamilton et al. proposed a random neighbor sampling strategy based on the GraphSAGE framework, which dynamically generates node embeddings by aggregating local neighborhood information [2]. Research based on data enhancement has been widely applied to this problem. For example, Su and Lee proposed the corpus-aware emotion cycle GAN (CAEmoCyGAN) method, whose core innovation lies in the introduction of the corpus-aware attention module, which enhances the separability of category boundaries by adjusting the distribution of generated data [3]; Luo et al. proposed a conditional Wasserstein GAN (c-WGAN) based method for electroencephalography (EEG) data enhancement, using gradient-penalized WGAN to generate labeled EEG data differential entropy (DE) [4]. In addition, for research based on loss functions, Li et al. proposed a gradient coordination mechanism (GHM), which uses the gradient density function to balance the weights of the model to distinguish between easy-to-distinguish and difficult-to-distinguish samples [5].

However, most of the existing reviews do not summarize the systematic comparison of data enhancement, loss function optimization, and other cross-technology routes to solve the data imbalance problem and lack a systematic summary of the differences in multi-scenario adaptation of text, speech, and EEG signals. For this reason, this paper analyzes seven representative methods (covering data enhancement, loss optimization, and multimodal generation) to systematically sort out the system of solutions to the data imbalance problem in the field of emotion recognition, compares the applicability boundaries of different methods in multimodal dialogues and EEG signals, and explores the optimization potentials of diffusion models, meta-learning, and other new technologies. The full paper is structured as follows: Section 2 Typical Technology Comparison and Analysis, section 3 Application Scenario Analysis, section 4 Challenges and Prospects, and section 5 Conclusions.

## 2. Typical technology comparison and analysis

### 2.1. The problem of category imbalance in emotion recognition datasets and its hazards

The common category imbalance in emotion recognition datasets is mainly manifested in the scarcity of a few category samples and semantic overlap between categories. For example, there are far fewer samples of "fear" and "disgust" than "neutral" and "happy" in the multimodal dialog emotion dataset. For example, there are far fewer samples of "fear" and "disgust" than of "neutral" and "happiness" in the multimodal dialog emotion dataset, and the boundary is further blurred by the partial overlap of semantics between different emotions such as "surprise" and "fear". The sample imbalance leads to a bias towards the majority class and ignores minority features, resulting in recognition bias and reduced generalization ability. Specifically, Poria et al. found that the F1 scores of the "disgust", "fear", and "sadness" categories were significantly lower than those of other categories in the MELD baseline experiments because of the extreme sample imbalance [6]. Zhang et al. also showed that "happy/neutral" samples accounted for the majority of facial expression data, while the rare "fearful/disgusted" samples were easily misclassified as positive emotions [7]. Li et al. further pointed out that ignoring this kind of category skew significantly reduces the generalization ability of the model, making it difficult to identify a few categories effectively [8].

## 2.2. Related solution introduction

This paper propose seven related state-of-the-art research techniques to address the problem of dataset imbalance in emotion recognition, and systematically sort out the latest technological advances in four dimensions: data optimization, loss function design, multimodal fusion, and complex scene adaptation, aiming to collaboratively explore the ideas to reduce the impact of imbalanced data on model performance through multi-dimensional collaboration.

### 2.2.1. Basic optimization at the data level

In the basic optimization at the data level, sample expansion and generation techniques are mainly used to improve the data set imbalance problem. Aiming at the characteristics of EEG signals with long time series and small sample size, the One-Dimensional Convolutional Neural Network and Bi-Directional Long Short Term Memory (1DCNN-BiLSTM) hybrid method proposed by K Singh et al adopts segmentation processing to expand the sample size, and combines the SMOTE (Synthetic minority oversampling technique) algorithm with the "Anger", "Sadness" and other A few classes are used for sample expansion to alleviate the class distribution bias while preserving the time-frequency characteristics of the original signal [9]. This method simplifies the traditional data enhancement process by extracting local features through 1D-CNN (one-dimensional convolutional neural network) and capturing global dependencies through bidirectional LSTM (long-short term memory neural networks), and achieves synergistic improvement of precision and recall in EEG emotion classification. In the scenario of multimodal emotion recognition in conversation (MMERC), the IMBA-MMERC (an effective framework designed to address the pervasive issue of class imbalance in MMERC) technique addresses the problem of semantic distortion caused by differences in session lengths, and proposes an improved SMOTE algorithm to achieve cross-modal sample synthesis: maximizing the constraints on the semantic consistency of text and speech modalities through mutual information, combining with contextual information to guide the generation process, so as to make the generated multimodal data more semantically consistent [8]; meanwhile, it introduces the classification goodness and encouragement loss to suppress the impact of data enhancement on the majority of the data. Loss is encouraged to suppress the negative impact of data enhancement on the majority class performance. This dynamic oversampling strategy effectively improves the generation quality of minority moods on datasets such as MELD, but the robustness to asynchronous modalities such as speech delay still needs to be optimized.

### 2.2.2. Model-level loss optimization

In terms of model loss optimization, dynamic weighting and contrast learning are used to enhance the recognition ability of minority classes.T Shi et al. proposed the dynamic weighting loss of the MultiEMO (a novel attention-based correlation-aware multimodal fusion framework) model by introducing the sample-weighted focal point contrast loss (SWFC) [10], which focuses on difficult-to-distinguish minority class samples by adjusting the temperature parameter $\tau$ and weight parameter $\alpha$, and significantly improves the discrimination between minority classes such as "fear" and semantically similar classes such as "anger - aversion" on both the IEMOCAP and the MELD datasets [11]. By adjusting the temperature parameter $\tau$ and weight parameter $\alpha$ to focus on the difficult-to-discriminate few categories, it can significantly improve the differentiation between the few categories such as "Fear" and the semantically similar categories, such as "Anger-Aversion" on the IEMOCAP and MELD datasets. The core advantage is to minimize category confusion through

inter-class distance, but the hyperparameter sensitivity and large-scale computational overhead are still outstanding issues. To address the representation limitations of graph networks in multimodal dialogues, W Ai et al [12], on the other hand, constructed dialog and event relation-aware graph convolutional neural network (DER-GCN) model graph comparison learning method design masked graph autoencoders (MGAEs), which improves the representation capability of graph convolutional neural network (GCN) and mitigates the problem of unbalanced class distribution by randomly masking and reconstructing the nodes and edges. It is also combined with a ternary loss function based on a balanced sampling strategy to optimize the class distribution imbalance problem while preserving the graph structure information as much as possible. The method achieves a weighted accuracy (WA) of 69.7% on the IEMOCAP dataset, but the high computational cost of self-supervised training with graph convolution restricts its deployment on edge devices [12].

### 2.2.3. Multimodal fusion and advanced generation

In addition to the basic improvement ideas, some studies have carried out certain innovations. To address the heterogeneity of multimodal data, Tao et al proposed the cross-modal generation technique that combines generative adversarial networks with cross-modal constraints (CBERL) to achieve semantically consistent sample generation, introduces the identity loss in the data augmentation method using generative adversarial networks (GANs), in order to ensure consistency in the distribution between the generated data and the original data [13]; and also introduces the sentiment classification loss and discrimination loss innovatively, guiding the generator to learn the complementary semantics of textual and visual and physiological signals. It also innovatively introduces affective classification loss and discrimination loss to induce the generator to learn the complementary semantics of textual, visual, and physiological signals, and enhances the GCN's ability to correct for the minority class of nodes through graph mask reconstruction. This scheme of combining adversarial generation and graph network improves the weighted F1 value by about 8% on the IEMOCAP dataset, which provides a new idea to solve the noise problem of multimodal data co-generation. In the field of EEG signals, Z Zhang et al proposed the Emotion Subspace Constrained Generative Adversarial Network (ESC-GAN) for the first time, which opens up a new way to use generative adversarial networks for EEG signal data augmentation and solve the problem of imbalance of emotion category data in EEG emotion recognition [14]; ESC-GAN pioneered the EEG signal cross-category editing technology, which converts the high-representation class signals into a few class samples by cross-category editing. Class signals into a few class samples, and design diversity-aware loss and boundary-aware loss to avoid pattern collapse. This method demonstrates on datasets such as DEAP [15] that data augmentation with ESC-GAN can help simple convolutional neural network classifiers achieve cutting-edge performance as well as generate samples for defense against adversarial attacks, and opens up new paths for solving the problem of EEG signal class imbalance.

### 2.2.4. Comprehensive programs for complex scenarios

Complex scenarios such as dynamic noise, multimodal asynchrony, and real-time interference pose higher challenges for emotion recognition. For example, physiological signals in public places are susceptible to environmental noise, and the temporal misalignment of multimodal data (e.g., speech, EEG) may destroy semantic associations. At this point, it is difficult for a single technique to meet the real-time and robustness requirements, and a comprehensive solution that integrates data generation, fusion modeling and hardware adaptation is needed. The ConvNeXt-Attention fusion

model (CNXAF) dynamic fusion framework proposed by A Li et al designs a ConvNeXt-Attention fusion model to process asynchronous data for multimodal physiological signals, combines with the conditional Self-Attention Generative Adversarial Network (c-SAGAN) to enhance signal characterization in noisy environments, and effectively copes with noise and real-time interference [16]. By dynamically adjusting the cross-modal attention weights, the scheme achieves effective suppression of noise interference on the DEAP dataset and the WESAD dataset [17], which provides a relocatable engineering idea for scenarios such as smart security and telemedicine. The core value of this kind of comprehensive scheme is to balance the algorithm complexity and practical deployment requirements, and to promote the technology from laboratory validation to real scenarios.

## 2.3. Performance comparison and evolution

Table 1: Comparison of key methods

| Method | Core methodology | Key Problems | Dataset | Performance |
|---|---|---|---|---|
| 1DCNN-BiLSTM | Data segmentation + SMOTE data enhancement | Limited data samples and inappropriate category distribution in EEG | DEAP | Directly through data augmentation, sample balancing improves data inequality |
| IMBA-MMERC | Improve SMOTE + maximize mutual information | Conversation Consistency Breaks and Most Class Performance Degradation | IEMOCAP、MELD、M3ED | Generated samples conform to dialog logic without semantic conflicts |
| MultiEMO | SWFC loss + cross-modal attention | Confusion of semantically similar categories (e.g., "sadness" vs. "frustration") | IEMOCAP、MELD | Loss function lightweight for real-time scenarios |
| DER-GCN | Self-supervised graph masking + Ternary loss function based on balanced sampling strategy | Inadequate modeling of long-tailed distributions and event relationships | MELD、IEMOCAP | Event Relationships Enhance Contextual Understanding |
| CBERL | GAN generation of multimodal data + graph mask optimization | Insufficient multimodal samples and class boundary blurring | MELD、IEMOCAP | High quality of cross-modal generation with clear inter-class boundaries |
| ESC-GAN | EEG signal cross-class transformation + boundary perception loss | EEG signal homogenization and fragile classification boundaries | DEAP、AMIGOS、SEED | Generated samples retain physiological signal timing characteristics |
| CNXAFF | c-SAGAN generates physiological signals + CNXAF fusion | Small sample overfitting and noise interference | DEAP、WESAD | Complex environment Anti-noise interference |

Table 1 shows that data imbalance solutions in the field of emotion recognition have gradually evolved from early single-point optimization (e.g., data oversampling) to multi-technology synergy (e.g., joint dynamic loss and generative compensation), and further evolved to scene adaptation (e.g., real-time feature compensation). This process reflects the paradigm shift from isolated patching to system balancing, and from static processing to dynamic adaptation, and provides the basis for hierarchical technology selection for medical, security, and other scenarios.

## 3. Application scenario analysis

### 3.1. Healthcare

In healthcare scenarios, emotion recognition can be used to monitor patients' psychological states (e.g., depression, pain, etc.), but clinical data tends to have a long-tailed distribution, and the scarcity of samples for minority classes of emotions (e.g., fear, aversion) makes it difficult for the model to discriminate. To address this problem, many approaches mitigate the imbalance through data augmentation and model optimization. The CBERL framework uses a multimodal GAN to generate minority class samples and is complemented by graph neural networks to optimize boundary learning, which improves the accuracy and F1 value of the model in the fear/aversion category by 10%-20% [13], and still has strong discriminative power even when clinical samples are missing. Enhancement methods such as ESC-GAN and CNXAF utilize conditional self-attentive GAN to synthesize emotion samples on EEG and other physiological signals, effectively expanding the data of minority classes; 1D-CNN-BiLSTM method obtains an average gain of 5.6%-34.88% by balancing the factor enhancement in a four-classification emotion recognition task [9], indicating good adaptability to the good adaptability to a small number of classes. In addition, multimodal fusion models such as MultiEMO significantly improve the recognition of difficult-to-distinguish minority categories by capturing the association between speech, text, and vision through cross-modal multi-head attention and combining with sample-weighted contrast loss [10], while the IMBA-MMERC framework validates effective correction of category bias on both English and Chinese multi-datasets through the sample generation and encouragement loss strategies [8], which can be applied to healthcare scenarios in different countries. Taken together, these methods show strong capabilities in generating minority category samples, optimizing boundary learning, and cross-modal information fusion, which are expected to significantly improve the accuracy of minority category emotion recognition in medical scenarios, and thus provide a solid technical guarantee for the clinical implementation of emotion monitoring systems.

### 3.2. Intelligent security

In smart security scenarios, multimodal emotion recognition can assist in monitoring abnormal emotions and behaviors, but the problem of category imbalance is further exacerbated by the complexity of environmental noise and extreme scarcity of extreme emotion samples. To address these challenges, DER-GCN introduces a graph-convolutional network for dialog and event relations, and combines it with comparative learning loss optimization, which significantly improves the average accuracy and F1 value on multimodal dialog datasets such as IEMOCAP and MELD [12], demonstrating that it still has good robustness in complex dialogs and noisy environments; however, DER-GCN is limited by the scarcity of samples from a few classes. However, DER-GCN is limited by the sparse samples of a few classes, and the misjudgment rate of the "surprise" class is still at a low level [12], whereas the CNXAF dynamic fusion framework achieves a stable accuracy of about 95% for both DEAP and WESAD when the enhancement factor is greater than 1, and it is expected to maintain a higher accuracy under different extreme environments by combining with the denoising technique of the CNXAF dynamic fusion framework [16]. In addition, the MultiEMO framework fuses speech, text, and visual information through a cross-modal bi-directional attention mechanism and utilizes sample-weighted contrast loss, which significantly improves the recognition of a few types of emotions [10], providing support for fine-grained emotion classification in security

scenarios. Taken together, these techniques enable the system to maintain high recognition performance in complex environments, showing potential for deployment in smart security.

## 4. Challenges and prospects

Despite the success of the above methods in improving data imbalance, the following core challenges remain in practical application:

First, existing methods suffer from the problem of modal alignment difficulties. In cross-modal dialog scenarios, the problem of timing asynchrony between text, speech and visual signals is extremely prominent. In reality, speech signals often have millisecond-level delays, which will undoubtedly lead to serious semantic misalignment. Existing methods such as CBERL and IMBA-MMERC, although able to reduce the alignment error to a certain extent, are not robust enough when facing large-scale asynchronous data. Therefore, the optimization strategy focuses on introducing a hierarchical cross-modal pre-training framework and embedding a temporal adaptive module in the Contrastive Language-Image Pre-training architecture, which can significantly enhance the performance of processing asynchronous data by dynamically adjusting the correlation between different modes in real-time based on the data characteristics through the self-attention mechanism.

Secondly, the insufficient quality of generated samples is also a difficult problem to solve. In the process of physiological signal processing, take EEG signal processing as an example, although ESC - GAN can reduce the $\delta$-wave power error, there is still a lot of room for improvement in terms of generation efficiency and compliance with physiological laws. The current generated signals generally suffer from noise interference and frequency band deviation. The key to improvement lies in combining the diffusion model with the domain expertise constraints, constructing a priori templates based on the clinical standards of physiological rhythms, and guiding the generation process; at the same time, introducing the adversarial training mechanism to suppress the noise interference, thus improving the signal-to-noise ratio of the signals, so that the generated signals can be more in line with the actual needs.

Meanwhile, dynamic boundary optimization is in urgent need. Traditional static loss functions are highly susceptible to performance degradation when migrating across scenarios. For example, MultiEMO adopts fixed class weights, and when the model migrates between different datasets, the F1 values of a few classes will drop significantly. While the SWFC loss function can dynamically adjust the weights, its hyperparameters rely on manual setting and cannot be fully automated for optimization. To solve this challenge, it is necessary to explore the dynamic adjustment framework of meta-learning, draw on the adaptive mechanism of Model-Agnostic Meta-Learning, update the parameters in real time based on the online confusion matrix, and introduce self-supervised comparative learning to mine the potential inter-class relationships, so as to improve the generalization ability of the model in different scenarios in an all-round way.

Finally, the high computational complexity seriously restricts the technology landing. In real-time scenario applications, the computational complexity of graph networks (DER-GCN) and multimodal fusion models is too high, which makes it difficult to be deployed on edge devices, and the effect that can be achieved by existing lightweight solutions is more limited. The feasible optimization path mainly contains two aspects: one is to use heterogeneous computing technology for acceleration, and select the most suitable hardware devices to execute for different computing tasks; the second is to use neural architecture search technology to design more compact models through automation, and strictly control the number of parameters of the model to meet the demand of real-time inference while ensuring the performance of the model.

## 5. Conclusion

This paper systematically analyzes seven types of representative methods for the data imbalance problem in emotion recognition, covering the technical routes of data enhancement, loss optimization, generative compensation and scene adaptation, revealing the evolution law from single-point patching to multilevel synergy, and proposing solution paths for the key challenges of generation quality optimization, dynamic boundary adjustment and lightweight deployment. It is found that the latest solution realizes a systematic breakthrough through data-model-generation linkage: at the data level, temporal signal reorganization and adversarial generation techniques significantly expand the sample size of a few classes; at the model level, dynamic loss function and graph comparison learning strengthen the ability to distinguish semantically ambiguous emotions; at the level of multimodal fusion and advanced generation, cross-modal semantic constraints are combined to generate adversarial networks with Emotion Subspace Constraint Generative Adversarial Network significantly improves the cross-scene generalization ability; while in some complex scenes, the dynamic fusion framework for multimodal physiological signal characteristics combined with generative techniques to enhance the signal characterization in noisy environments to effectively cope with noise and real-time interference. In practical application scenarios (e.g., healthcare, smart security), different technologies can be combined to achieve improved accuracy and stable performance in emotion recognition. Although the challenges of practical applications still exist, the imbalance problem of emotion recognition datasets is gradually and systematically overcome with the continuous and synergistic development of the technology. Future research needs to further integrate the high-fidelity generation capability of diffusion modeling with the scene-adaptive properties of meta-learning, and at the same time establish multimodal alignment standards and medical ethical review mechanisms, in order to promote the reliable landing of emotion recognition technology in medical monitoring, intelligent security, and other scenarios.

## References

[1]  J. Deng and F. Ren, "A Survey of Textual Emotion Recognition and Its Challenges, " in IEEE Transactions on Affective Computing, vol. 14, no. 1, pp. 49-67, 1 Jan.-March 2023, doi: 10.1109/TAFFC.2021.3053275.

[2]  W. Hamilton, Z. Ying and J. Leskovec, "Inductive representation learning on large graphs" in Proc. Adv. Neural Inf. Process. Syst., MIT Press, vol. 30, 2017.

[3]  B.-H. Su and C.-C. Lee, "Unsupervised cross-corpus speech emotion recognition using a multi-source cycle-GAN", IEEE Trans. Affect. Comput., vol. 14, no. 3, pp. 1991-2004, Jul./Sep. 2023, [online] Available: .

[4]  Y. Luo and B.-L. Lu, "Eeg data augmentation for emotion recognition using a conditional wasserstein gan, " in 2018 40th annual international conference of the IEEE engineering in medicine and biology society (EMBC). IEEE, 2018, pp. 2535–2538.

[5]  B. Li, Y. Liu and X. Wang, "Gradient harmonized single-stage detector", Proc. AAAI Conf. Artif. Intell., vol. 33, no. 01, pp. 8577-8584, 2019.

[6]  S. Poria, D. Hazarika, N. Majumder, G. Naik, E. Cambria and R. Mihalcea, "MELD: A multimodal multi-party dataset for emotion recognition in conversations", Proc. 57th Annu. Meeting Assoc. Comput. Linguistics, pp. 527-536, 2019.

[7]  Zhang Y, Li Y, Liu X, et al. Leave no stone unturned: Mine extra knowledge for imbalanced facial expression recognition [J]. Advances in Neural Information Processing Systems, 2023, 36: 14414-14426.

[8]  Li Q, Huang P, Xu Y, et al. Generating and encouraging: An effective framework for solving class imbalance in multimodal emotion recognition conversation [J]. Engineering Applications of Artificial Intelligence, 2024, 133: 108523.

[9]  Singh K, Ahirwal M K, Pandey M. Subject wise data augmentation based on balancing factor for quaternary emotion recognition through hybrid deep learning model [J]. Biomedical Signal Processing and Control, 2023, 86: 105075.

[10] Shi T, Huang S L. MultiEMO: An attention-based correlation-aware multimodal fusion framework for emotion recognition in conversations [C]//Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). 2023: 14752-14766.

[11] C. Busso et al., "LEMOCAP: Interactive emotional dyadic motion capture database", Lang. Resour. Eval., vol. 42, no. 4, pp. 335-359, 2008, [online] Available: .

[12] Ai W, Shou Y, Meng T, et al. Der-gcn: Dialog and event relation-aware graph convolutional neural network for multimodal dialog emotion recognition [J]. IEEE Transactions on Neural Networks and Learning Systems, 2024.

[13] Meng T, Shou Y, Ai W, et al. Deep imbalanced learning for multimodal emotion recognition in conversations [J]. IEEE Transactions on Artificial Intelligence, 2024.

[14] Zhang Z, Zhong S, Liu Y. Beyond mimicking under-represented emotions: deep data augmentation with emotional subspace constraints for EEG-based emotion recognition [C]//Proceedings of the AAAI conference on artificial intelligence. 2024, 38(9): 10252-10260.

[15] S. Koelstra, C. Muhl, M. Soleymani, J.-S. Lee, A. Yazdani, T. Ebrahimi, T. Pun, A. Nijholt, and I. Patras, "Deap: A database for emotion analysis; using physiological signals, " IEEE transactions on affective computing, vol. 3, no. 1, pp. 18–31, 2011.

[16] Li A, Wu M, Ouyang R, et al. A Multimodal-Driven Fusion Data Augmentation Framework for Emotion Recognition [J]. IEEE Transactions on Artificial Intelligence, 2025.

[17] P. Schmidt, A. Reiss, R. Duerichen, C. Marberger, and K. Van Laerhoven, "Introducing wesad, a multimodal dataset for wearable stress and affect detection, " in Proceedings of the 20th ACM international conference on multimodal interaction, 2018, pp. 400–408.