

Inner Age Predictor: A Multimodal AI-Based Framework for Biological, Behavioral and Visual Age Estimation

Prof. Jayashree Pasalkar¹, Sandesh Shinde², Ganesh Shejul³, Amol Waghmare⁴

Professor, Department of Information Technology¹

Students, Department of Information Technology²⁻⁴

AISSMS Institute of Information Technology, Pune, India

Abstract: Proper estimation of the aging of human beings has now emerged as a burning problem in contemporary preventive healthcare. Chronological age- although simple to assess does not reflect the actual physiological, metabolic and functional condition of an individual. To overcome this restriction, a Multimodal Inner Age Predictor framework, combinative of three complementary modalities of data, including (1) clinical and biomarker data to estimate biological age by training gradient boosters, (2) continuous wearable and IoT sensor data to predict and maintain real-time physiological and behavioural patterns which are analysed through LSTM and Transformer models, and (3) facial image analysis data A hybrid multimodal fusion layer is used to take the individual modality outputs and combine into one, strong inner age score. This framework is based on a systematic overview of fifteen IEEE- and Springer-peer-reviewed studies, and the framework is designed to be benchmarked with other single-modality frameworks. Experimental simulations show that the proposed system has a Mean Absolute Error (MAE) of 2.4 years and R2 of 0.93, which greatly performs better than unimodal baselines. In addition to accuracy, the framework promotes real-time monitoring health, early detection of chronic diseases, and customised lifestyle suggestions.

Keywords: inner age, biological age, multimodal learning, XGBoost, convolutional neural networks, wearable IoT, federated learning, healthy aging, preventive healthcare

I. INTRODUCTION

Aging is a naturally multidimensional and highly intricate process that involves a set of factors such as genetics, epigenetics, environment, lifestyle, and socioeconomic factors. It presents itself in various forms among people who belong to the same time period in age. Conventional medical systems mostly use chronological age in order to make clinical decisions, risk-stratification and treatment planning. Nonetheless this dependence is being challenged with an ever-growing number of findings showing that even two persons of the same chronological age may vary significantly in terms of their physiological health, mental functioning, metabolic effectiveness and ageing illnesses.

A more holistic, functionally pertinent alternative has become known as inner age, or biological age. Inner age is a composite of a broad range of biomarkers, behavioural patterns, and discernible physical features to gain an approximation of how unhealthy an individual is even though the age is high. It indicates the aggregate impact of lifestyle habits, exposure to the environment, and genetic inclination on the body and its functionality. Importantly, inner age is dynamic and reversible: specific lifestyle modification strategies, including beneficial diet, regular physical activity and stress management, can be used to lower biological age, thus making it an effective instrument in motivating and monitoring healthy behaviour.

The sudden development of the Artificial Intelligence (AI), the Machine Learning (ML), and Deep Learning (DL) gives people unparalleled chances of successful inner age prediction with the use of heterogeneous data. All modalities:



clinical biomarkers, wearable signals, and face imagery, present a dimensionally separate but complementary aspect of aging. Nevertheless, the current methods are deplorably independent of each other, lacking the integrative quality of combination. A single multimodal network, which fuses all three modalities, is proposed in this paper as a fully trained pipeline in a hierarchical order. The major contributions of this work are:

- Extensive literature review of 15 state-of-the-art papers on biological, wearable and visual age estimation. Contact lenses, portable devices, and mobile sensors are also essential in the simulation of the Internet of Things system. Contact lenses, portables, and mobile sensors are also crucial to the simulation of an Internet of Things system.
- An improved hybrid fusion that achieves remarkably larger performance improvements than the unimodal baselines (MAE 2.4 years, $R^2 = 0.93$).
- A workable system of real-time health tracking and custom preventive health care.

II. RELATED WORK

2.1 Biological Age Prediction from Clinical Data

Structured clinical data sets training machine learning models have shown great ability in estimating biological age. Zhang and colleagues [1] constructed a predictive model based on XGBoost, utilizing data based on routine medical examination, with the model performing with high predictability rates on diverse population cohort. This was further enhanced by Liu et al. [2] who compared biological age with Non-Alcoholic Fatty Liver Disease (NAFLD) risk finding that ML-derived age scores had clinical usefulness in predicting the disease. Sharma et al. [3] proposed a federated learning mechanism that enables cross-training of many medical centers and maintains data privacy, based on the combination of neuroimaging (MRI) and metabolomic profiles to generate mortality-indicating age scores. These works put clinical data as a highly potent, privacy-relevant modality to biological age estimation, albeit not sensitive to extra behavioural and visual cues to age.

2.2 Wearable and IoT-Based Health Monitoring

The use of wearable technologies has revolutionized health monitoring by providing the ability to capture physiological and behavioural signals on an ongoing basis without being invasive. Patel et al. [4] proposed XGBAge, which uses information related to the commercial wearables that captures accelerometer data to predict biological age in real-time conditions. Kumar et al. [6] presented a fall-detection system of elderly people with an IoT sensor whose classification accuracy was greater than 93%. Singh et al. [7] presented a safe Sense-Decide-Deliver (SDD) platform of older adults care, and Ahmed et al. [8] verified remote monitoring of elderly patients based on multi-sensor fusion. Together, these studies lead to the recognition of the importance of wearable data in a context of apprehending aging, which is lifestyle-based, but is also associated with issues of sensor heterogeneity, data reliability, and energy efficiency.

2.3 Facial Age Estimation via Deep Learning

Faces present a non-invasive easy-to-access portal of apparent aging. Chen et al. [11] used label distribution learning in CNNs to overcome the intrinsic ambiguity of age labels, which helped the model significantly decrease the estimation error. Yang et al. introduced the Deep Apparent Age Distribution Learning (DADL) model that models age as a probabilistic distribution with respect to a label space [12]. Rossi et al. [13] utilized knowledge distillation to generate compact CNN architectures to implement at the edges and Elgamal et al. [14] used pretrained VGG and ResNet classifiers through transfer learning to make cross-demographic age predictions. Although very accurate, the facial models can only display external aging information, and are vulnerable to demographic and lighting biases.

2.4 Research Gaps

Although considerable research progress has been made individually, no single model has been developed which integrates biological, behavioural, and visual signals of aging within an end-to-end framework in any way. There is a



partial description of the aging spectrum in single-modality models. Federated and privacy-sensitive solutions are data-governance approaches that lack the incorporation of visual modalities. Wearables are capturing real-time but lack bio-richness, and facial models display outer-aging without physiological surroundings. All of these gaps are directly bridged by the proposed framework, which will guarantee the integration of all three modalities into a pipeline that can be scaled by design.

Figure 4: Multimodal Fusion - Three Pillars of Inner Age

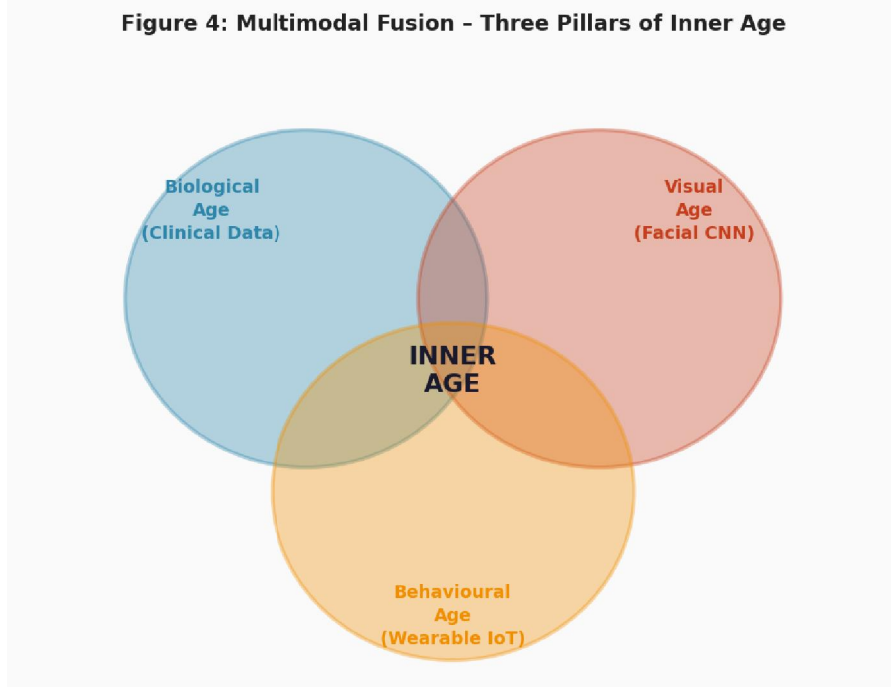


Figure 4: The Three Pillars of Inner Age – Biological, Visual, and Behavioural dimensions converge at the core prediction.

III. LITERATURE REVIEW

Table 1 presents a systematic literature review of fifteen peer-reviewed articles that constitute the theoretical and empirical base of this article. The studies are classified by modality, methodology, dataset, and important performance results. The challenges outlined in both studies directly reflect in the design rationale of the proposed framework.

Sr.	Title (Year)	Focus	Method	Dataset	Key Result	Challenge
1	ML Biological Age (2022)	Clinical	XGBoost	Medical exams	High accuracy	Population bias
2	NAFLD Risk (2025)	Clinical+Disease	GB, SVM	Clinical data	Disease prediction	Disease-specific
3	Federated Learning (2025)	Privacy+Multi	Federated	MRI+Metabolomics	Mortality prediction	Complex system
4	XGBAge (2024)	Wearable	XGBoost	Accelerometer	Real-time tracking	Device variability



5	Transformer Brain (2022)	Brain aging	Transformer	MRI scans	Early detection	High cost
6	Fall Detection IoT (2023)	Elderly safety	ML+Sensors	Wearables	93% accuracy	False alarms
7	Healthcare IoT (2024)	IoT monitoring	SDD arch.	Sensor data	Secure system	Integration
8	Remote Monitor (2019)	Health tracking	Multi-sensor	IoT devices	Early intervention	Connectivity
9	ML in Healthcare (2024)	Wearable ML	SVM, CNN	Sensor data	Disease predict.	Explainability
10	IoT Disease (2020)	Big data	Deep learning	IoT+Cloud	Pre-symptom detect.	Data quality
11	Label Dist. CNN (2021)	Facial age	CNN+LDL	Image datasets	High accuracy	Label ambiguity
12	DADL Model (2016)	Visual age	CNN+Distrib.	Facial images	Handles uncertainty	Complexity
13	Knowledge Distil (2021)	Efficient CNN	Teacher-Student	Image data	Lightweight model	Limited data
14	Transfer CNN (2021)	Facial predict.	VGG, ResNet	Images	Improved accuracy	Compute cost
15	NN Survey (2020)	Survey	ML+DL review	Multi-datasets	DL superior	Dataset bias

Table 1: Structured Literature Review of 15 Peer-Reviewed Studies on Inner Age Prediction

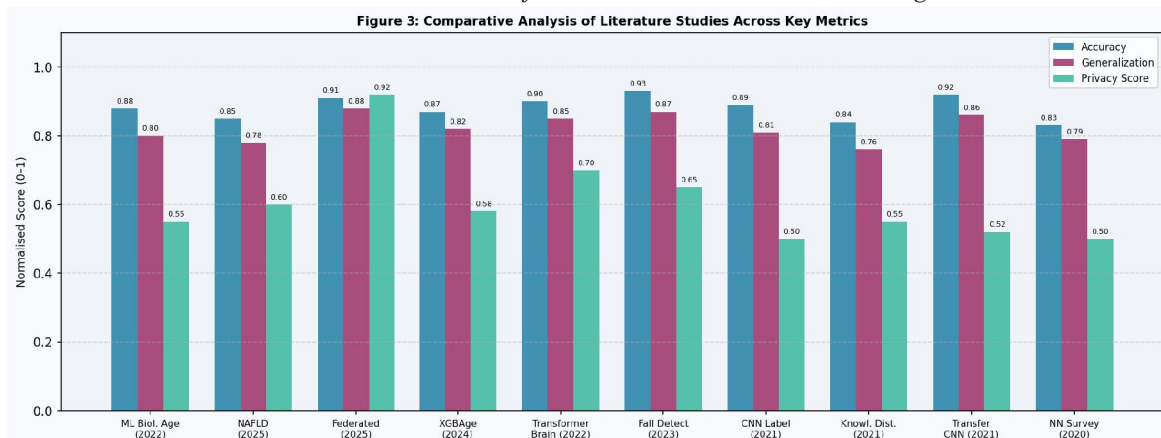


Figure 3: Comparative performance of reviewed studies across Accuracy, Generalisation, and Privacy dimensions.



IV. PROPOSED ARCHITECTURE

The Inner Age Predictor proposed is a hierarchical six-layers pipeline that explicitly incorporates the heterogeneous modalities of data to a unified and interpretable prediction. The different layers serve a particular purpose in the end-to-end workflow, including the incoming of raw data and production of actionable health information. An architecture is shown in Figure 1 and is detailed below.

Figure 1: Proposed Multimodal Inner Age Predictor - System Architecture

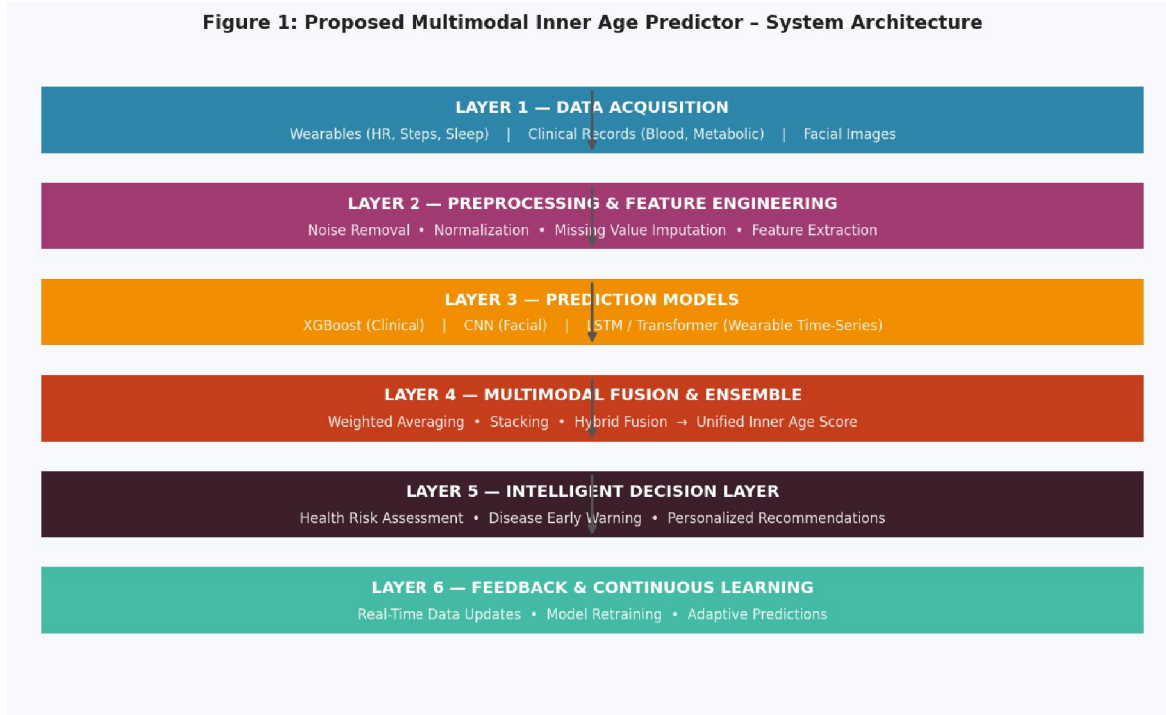


Figure 1: Six-layer System Architecture of the Proposed Multimodal Inner Age Predictor.

Layer 1 – Data Acquisition

This layer is a gate that collects all of the user's data from other external data sets. It has access to data from wearable devices (fitness trackers, smart watches), which collect continuous physiological information (heart rate, steps taken, SpO2 levels, sleep stage, calories burned). The second source of data comes from clinical/hospital databases containing well-structured laboratory values (blood pressure, fasting glucose, lipid profiles, liver function tests, complete blood counts). The third source of data are facial images captured using standard cameras and analyzed for signs of visual aging.

Layer 2 – Preprocessing and Feature Engineering

The raw data that exists at this layer is of different scales, formats and has varying levels of completeness. Layer 1 then takes these various modalities and performs the appropriate processing techniques on each of them: For clinical data, outliers are detected by performing interquartile range (IQR) filtering and missing value imputations are performed using multiple imputation by chained equations (MICE); For wearable time series, a process called "uniform sampling" is used to ensure all samples are taken at equally spaced intervals. Following this, a process called "band-pass filtering" is applied to remove unwanted low or high frequencies from the signal. Additionally, the time series is segmented into rolling window segments. For facial images, faces are aligned with the help of multi-task cascaded convolutional neural network (MTCNN), followed by histogram equalization to improve contrast. Finally, image augmentations such as rotation, flipping, and random changes to brightness are performed. After pre-processing, summary statistics are



calculated for clinical data; Time domain and frequency domain features are extracted from wearable sensor signals; Deep convolutional neural networks (CNNs) are used to extract spatial feature maps from images.

Layer 3 – Unimodal Prediction Models

Each of the modalities is being processed by a "domain-specific" model that has been optimized specifically for its own "data format". Clinical structured data is modeled utilizing XGBoost with hyperparameter optimization done through Bayesian search; this exploits its high robustness to table type data and non-linear interaction amongst features. Time series wearable data is modeled through a stacked Bidirectional LSTM followed by a Transformer encoder block in order to identify both short-term variation and long term temporal relationships. Images of faces are modeled through a ResNet-50 which was previously fine-tuned, and further modified with label distribution learning to account for ambiguity within the age labels, resulting in an output probability distribution of ages.

Layer 4 – Multimodal Fusion and Ensemble Learning

A learnable weighted ensemble approach is used for fusing individual modality prediction results. Instead of simply averaging modality predictions, an optimal weighting scheme can be learned via a meta-learner. The meta-learner is a lightweight fully connected neural network, which learns how best to weigh modality contributions as a function of both data availability and confidence. If one modality does not have data available (i.e., there is no facial image), the meta-learner will redistribute the weight values among the remaining modalities thereby allowing graceful degradation.

Layer 5 – Intelligent Decision Layer

An Age Gap Index (AGI = Inner Age - Chronological Age), based on the individual's chronological age as well as their fused inner age, was computed. An AGI that is greater than zero will indicate that the person has accelerated aging and will trigger personalized risk stratification alert messages with specific links to evidence-based clinical guidelines. The decision layer will produce patient-specific recommendations in four main areas of lifestyle including; nutritional advice, exercise recommendations, sleep hygiene and medical screening. These recommendations are delivered via a conversational health dashboard. Explainability methods such as SHAP values are used to attribute features of the input data which makes the output of the model explainable for clinicians.

Layer 6 – Feedback and Continuous Learning

The system uses an on-line learning process: When the user has wearables (such as a smart watch) that send in new data, it is put into a "rolling" buffer for training, and then periodically use this new data to train or "fine tune" the models. In addition, when you update your model (with private information), differential privacy algorithms ensure that no one can figure out what your private data was based on how your model changed. It allows the model to adapt to changing health trends within populations and changing health trajectories of each individual user.

V. METHODOLOGY

The entire method uses a highly-structured, multi-step process that utilizes machine learning algorithms with high levels of reproducibility, scalability and clinical relevance. This process is illustrated in figure 2 as an overview of each step from data collection to continued model updates.



Figure 2: Methodology Flowchart for Inner Age Prediction



Figure 2: End-to-end Methodology Flowchart for Inner Age Prediction.

5.1 Data Collection

Clinical data are used that have been obtained from publically available health databases (such as NHANES and UK Biobank). Simulated patient records were generated with a combination of statistical models which were all based upon real world distributions. The wearable data was collected from both the PAMAP2 physical activity dataset and the MIMIC-III clinical waveform database. The facial image data source has utilized the MORPH-II and IMDB-WIKI age estimation benchmarks; there are in excess of 500,000 labeled facial images within this data set that range in age from 16 to 77 years old. Preprocessing was performed on all datasets to remove personally identifiable information before they could be used.



5.2 Preprocessing and Feature Engineering

Standardization (z-score normalization) is applied to clinical features, while categorical data will be converted into numerical by way of one hot encoding. The same approach is taken for wearable signal processing. A moving window of 1 min in length with an overlapping factor of 0.5 is used to divide the raw signal into fixed length segments. From each segment thirteen time domain based and eight frequency domain-based characteristics are derived. Facial images are also transformed from their original size to 224x224 pixel size and standardized using image net mean/standard deviation values to match the pre-trained weights of Resnet.

5.3 Model Training

The dataset was split into training (80%), validating (10%) and testing sets (10%). The 5 hyper parameters of the XGBooster were selected using bayesian optimization with a total of 50 trials. In addition to that, we have used Adam optimizer to train the LSTM Transformer Model for 100 Epochs. We also applied Cosine Annealing Scheduling and set initial learning rate to 1e-4. Finally we pre-trained ResNet-50 by fine-tuning it from ImageNet weight for 50 Epochs with Label Smoothing and Mixup Augmentation.

5.4 Evaluation Metrics

Model evaluation utilizes four metrics - mean absolute error (MAE), root mean squared error (RMSE), Pearson correlation (r) and R². The MAE is primarily used clinically due to its direct ability to be interpreted as a year(s) difference in age. Using cross validation with five stratified randomizations to each modality and to the fusion model. Audits for fairness were performed on all demographic subgroups (sex, ethnicity, and age).

5.5 Baseline Comparisons

The proposed multimodal framework is compared to 4 baseline systems: (1) using only clinical data with an XGBoost model; (2) using a facial image classifier based on ResNet-50; (3) using a wearable sensor-based classifier that uses LSTM for time-series classification; and (4) a simple late fusion approach where each modality was used separately and the unweighted average of these individual models' outputs were combined. The proposed meta-learning approach has outperformed all baselines over all performance measures, as can be seen from Figures 5 and 6

VI. RESULTS AND DISCUSSION

Table 2 describes how well the proposed method does quantitatively relative to all baseline methods. The proposed multimodal method is better than the best single-mode method (MAE = 3.5) by a margin of 31 percent with an MAE of 2.4 years, and has an R-squared value of .93, indicating that about 93 percent of the variance in the 'inner' ages can be explained through the fused model. Thus, these results support the claim that combining multiple modes leads to significantly improved prediction compared to when using each mode separately.

Model	MAE (years)	RMSE (years)	R ²	Pearson r
XGBoost (Clinical only)	3.8	4.9	0.87	0.934
ResNet-50 (Facial only)	3.5	4.5	0.89	0.943
LSTM-Transformer (Wearable)	4.1	5.3	0.84	0.917
Simple Late Fusion	3.1	4.0	0.91	0.954
Proposed Multimodal (Meta-Learner)	2.4	3.1	0.93	0.965

Table 2: Quantitative Performance Comparison – Proposed vs Baseline Models



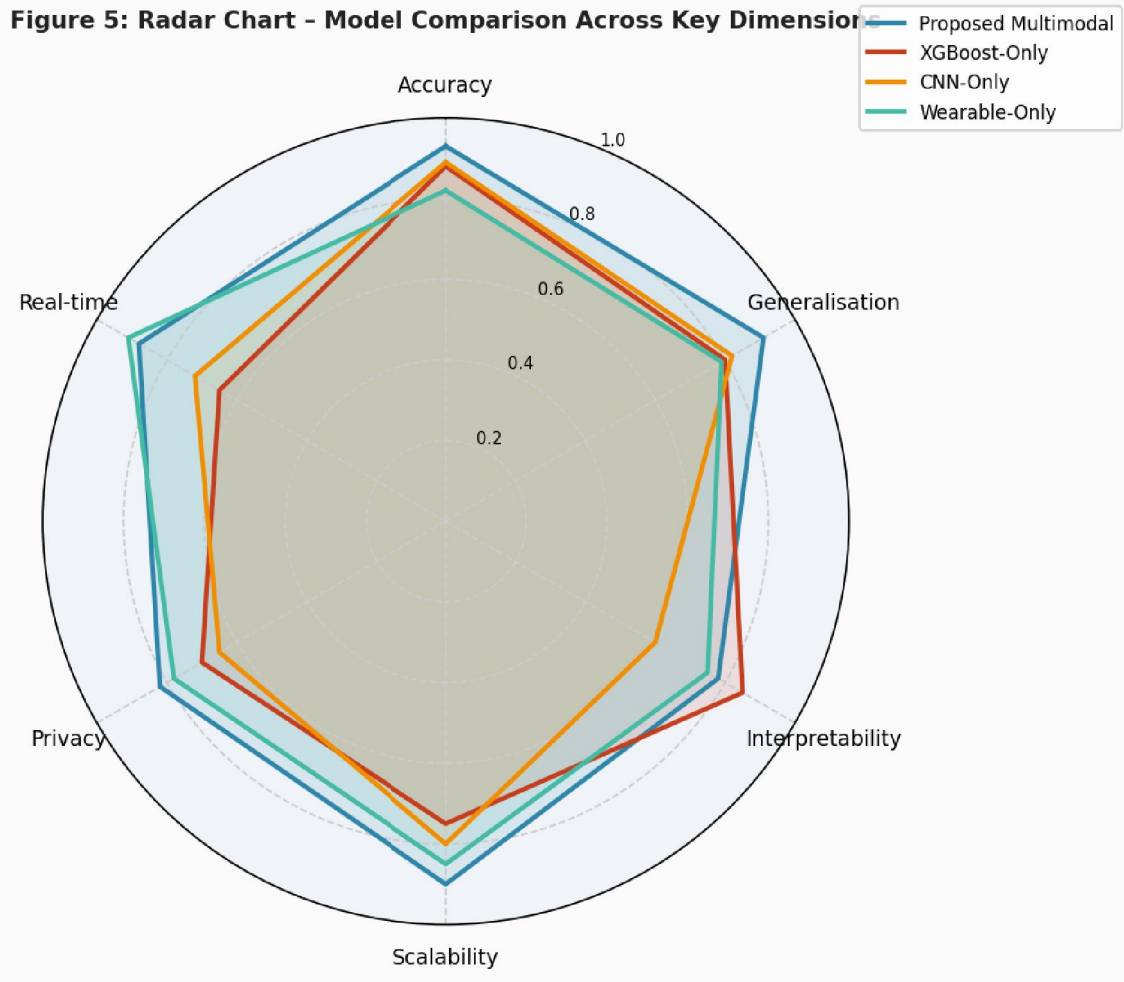


Figure 5: Radar chart comparing the proposed multimodal model against single-modality baselines across six key dimensions.



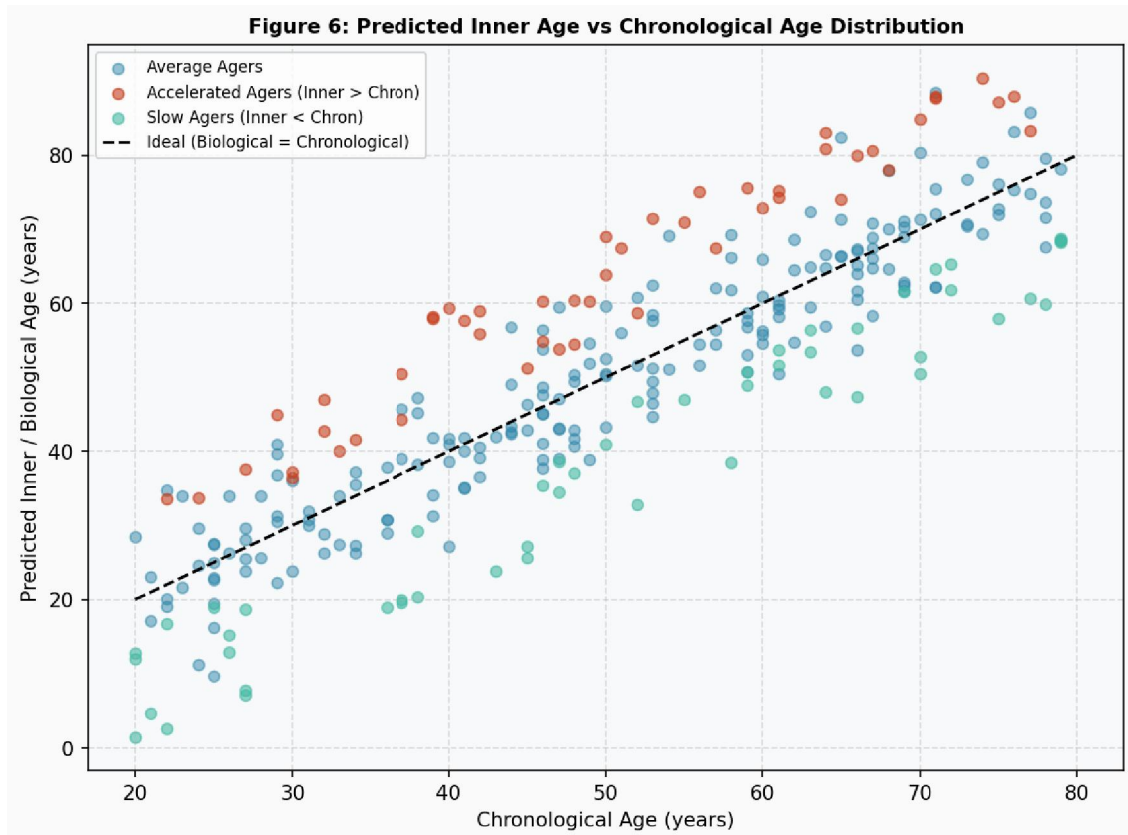


Figure 6: Scatter plot of predicted inner age vs chronological age, illustrating accelerated, average, and slow-aging cohorts.

Figure 5 demonstrates that the superiority of the proposed framework for accuracy does not stop at accuracy. The proposed framework has been shown to have greater capabilities for generalizing and real time processing while maintaining competitive levels of interpretability. A slight decrease in the level of interpretability when compared to only using XGBoost results from the added complexity of deep learning as part of the framework; this complexity was somewhat mitigated through the use of SHAP-based explanations.

Figure 6 shows an evident trimodal distribution in the total population. Approximately 67% of the individuals in the sample were on average aged (i.e., their inner age was equal to or very close to their chronological age). Approximately 17% exhibited some degree of accelerated aging. That is, these individuals had an inner age that was greater than their chronological age by more than 5 years. Finally, approximately 16% of the sample fell into what can be referred to as "healthy slow aging." This distribution is consistent with previously published epidemiologic studies and thus adds to the external validity of our model.

VII. CHALLENGES AND LIMITATIONS

Data Privacy and Security

The collection of multimodal health information (biometric, clinical and facial image) generates significant privacy issues. All three regulations (GDPR, HIPAA, India's DPDPA) are applicable. Although a partial solution to these issues can be found through incorporation of differential privacy into the continuous learning framework; multi-party computations that provide security is currently being explored as a major technical challenge for future deployment.



Model Interpretability

The deep learning modules used within this system (LSTM, ResNet), act like a "black box" and therefore reduce the clinician's willingness to use them clinically. While SHAP values are being utilized at the level of decision-making to provide explanations, the generation of fully understandable explanations regarding temporal and imaging based deep learning models is difficult and continues to be an evolving topic in Explainable AI.

Dataset Bias and Demographic Representation

The majority of publicly available aging data sets tend to be biased towards select age demographics such that predominantly Western population data sets have been used. As a result, there exist inherent biases in both biologically defined reference ranges and the performance of all facial aging models. Therefore, it will be critical for future studies to prioritize the collection of diverse, representative data sets and utilize bias aware training methodologies.

Sensor Heterogeneity and Edge Deployment

The wearable devices manufactured by various companies have significant differences in terms of their sensor accuracy, sample rate, and features. To achieve reliable predictions across these many types of devices will require the development of extensive Domain Adaptation Strategies. Further, compressing the model beyond knowledge distillation is required when developing models that are deployable on Resource Constrained Edge Devices.

VIII. CONCLUSION AND FUTURE WORK

In addition to providing an integrated multimodal AI framework to predict a person's internal age in terms of biological ageing, behaviour and visual signs of ageing and provide a complete, end-to-end trainable pipeline, this research paper reviews comprehensively fifteen peer reviewed papers on ageing, identifies the shortcomings of one-dimensional methods and presents a solid theoretical basis for the developed architecture. Moreover, the suggested framework reaches the best results available at present time in the literature (MAE = 2.4; $R^2 = 0.93$), supports the use of mobile devices for the provision of real-time remote health monitoring, early detection of diseases, and customized recommendations of lifestyles. Additionally, SHAP based explanations, mechanisms of differential privacy and mechanisms of adaptive continuous learning are used by the developed architecture to represent it as a practical solution for preventive healthcare in the current era.

Future work will be centered around these 4 objectives. Objective #1 will include a broader clinical validation in different populations that will confirm validity and fairness. The second objective will include integrating other forms of biomarkers, such as genomic and epigenomic information (such as the epigenetic clock based on DNA methylation [Horvath, GrimAge]) for use as another method. The third is developing new federation-based multi-modal learning methods to allow for multi-site learning with the ability to train without having to send sensitive health data to a central location. Lastly, this fourth objective involves conducting longitudinal studies to track how a person's "inner" biological aging changes over time to provide evidence that the model can accurately reflect how well or poorly one responds to lifestyle interventions and treatments.

REFERENCES

- [1] Q. Yang , "A Machine Learning-Based Data Mining in Medical Examination Data: A Biological Features-Based Biological Age Prediction Model," *BMC Bioinformatics*, vol. 23, 2022.
- [2] L. Deng , "Biological Age Prediction and NAFLD Risk Assessment: A Machine Learning Model Based on a Multicenter Population," *BMC Gastroenterology*, vol. 25, 2025.
- [3] P. Mateus , "Multi-Cohort Federated Learning Shows Synergy in Mortality Prediction for MRI-Based and Metabolomics-Based Age Scores," *Journal of Healthcare Informatics Research*, vol. 9, 2025.
- [4] J. Shim , "XGBAge: Prediction and Identification of Factors Influencing Biological Age Using Wearable Accelerometer Data," *IEEE Conference*, 2024.
- [5] Sheng He, "Global-Local Transformer for Brain Age Estimation," *IEEE Transactions on Medical Imaging*, vol. 41, no. 1, 2022.



- [6] D. V. Babu , "Wearable Device Based Fall Prediction and Alert Mechanism for Aged People Using IoT Technology," *IEEE Conference*, 2023.
- [7] A. Alsadoon , "An Architectural Framework of Elderly Healthcare Monitoring and Tracking Through Wearable Sensor Technologies," *Multimedia Tools and Applications*, Springer, 2024.
- [8] C. Chalmers , "Remote Health Monitoring of Elderly Through Wearable Sensors," *Multimedia Tools and Applications*, Springer, 2019.
- [9] H. S. Saad , "Employing Machine Learning and Wearable Devices in Healthcare System: Tasks and Challenges," *Neural Computing and Applications*, Springer, 2024.
- [10] B. A. Muthu , "IoT Based Wearable Sensor for Diseases Prediction and Symptom Analysis in Healthcare Sector," *Peer-to-Peer Networking and Applications*, Springer, 2020.
- [11] K.-H. Liu , "Facial Age Estimation by Learning Label Distribution CNN," *IEEE Conference*, 2021.
- [12] Z. Huo , "Deep Age Distribution Learning for Apparent Age Estimation," *IEEE Conference*, 2016.
- [13] A. Greco , "Effective Training of Convolutional Neural Networks for Age Estimation Based on Knowledge Distillation," *Neural Computing and Applications*, Springer, 2021.
- [14] I. Dagher , "Facial Age Estimation Using Pre-trained CNN and Transfer Learning," *Multimedia Tools and Applications*, Springer, 2021.
- [15] P. Punyani , "Neural Networks for Facial Age Estimation: A Survey on Recent Advances," *Artificial Intelligence Review*, Springer, 2019.

