

# Age Engagement Estimation in Human-Centered AI

Rajni Dubey<sup>1</sup> and H N Verma<sup>2</sup>

M.Tech.(CSE) Student, Dept. of CSE<sup>1</sup>

Associate Professor, Dept. of CSE<sup>2</sup>

ITM University, Gwalior, Madhya Pradesh, India

beti1cs21002@itm university.ac.in and hnverma.mca@itm university.ac.in

ORCID: 0000-0003-3268-0787

**Abstract:** *Human-Centered AI (HCAI) demands models that are not only accurate but also reliable, fair and intelligible to the people they serve. This paper tackles a joint problem in term age–engagement estimation: simultaneously inferring a user’s apparent age and moment-to-moment engagement state from multi-modal signals (face, gaze, posture, interaction logs) so that interfaces can adapt responsibly. It motivate the task in safety-, education- and accessibility-critical settings, survey post-2020 advances in facial age estimation and automatic engagement analysis and propose a unified, privacy-aware learning objective with fairness regularization. Our method integrates ordinal age modeling, temporal engagement inference and human- centered constraints (documentation, transparency, controllability). On benchmark datasets and a controlled pilot, the approach produces competitive age MAE and robust engagement F1 while reducing disparity across age groups. It presents ablations and qualitative analyses that relate attention maps and decision rules to meaningful behavioral cues. It concludes with a roadmap for deploying age–engagement systems that meet HCAI standards in the wild*

**Keywords:** Age estimation; Engagement recognition; Human- Centered AI; Fairness; Transparency; Ordinal regression; Multimodal learning; Privacy; Responsible AI; Interpretability

## I. INTRODUCTION

Human-Centered AI reframes “accuracy-at-all-costs” machine learning by foregrounding human goals—safety, dignity and agency—alongside performance. In this framing, study age–engagement estimation: jointly predicting a person’s apparent facial age (useful for content gating, accessibility defaults) and their engagement (useful for timing prompts, pacing content). The pairing matters because naïvely optimizing one can harm the other (e.g., high- contrast animations that raise engagement but disadvantage older users with visual strain). Contemporary HCAI work emphasizes reliability, safety and user control; It extend that ethos to model design, evaluation and reporting. Recent literature shows surging interest in engagement detection for learning and HCI, with computer-vision and multimodal pipelines outperforming manual checklists and a parallel stream in facial age estimation that blends ordinal and generative paradigms. These trajectories motivate a single, regularized objective that balances utility with constraints on fairness and privacy. Shneiderman’s HCAI program; systematic engagement reviews; recent age-estimation advances [1]. It formalizes the problem as follows. Let  $x_{vis}$  be visual frames,  $x_{int}$  interaction logs (keystrokes/mouse) and  $y_{age} \in \{0, \dots, 100\}$  an apparent-age label (or distribution), with engagement  $y_{eng} \in \{0, 1, 2\}$  (disengaged/neutral/engaged). A shared encoder produces representations  $h$ . It minimizes a joint risk:

$$\mathcal{L} = \lambda_{age} \mathcal{L}_{ord}(y_{age}, \hat{y}^{age}) + \lambda_{eng} \mathcal{L}_{ce}(y_{eng}, \hat{y}^{eng}) + \lambda_{fair} \Omega(\hat{y}^{age}; g) + \lambda_{priv} \Psi(h), \quad (1)$$

ordinal age loss engagement cross-entropy group fairness privacy

where  $\mathcal{L}_{ord}$  is an ordinal regression loss for age,  $\mathcal{L}_{ce}$  is a classification loss for engagement,  $\Omega$  penalizes disparities across protected groups  $g$  (e.g., apparent-age bins) and  $\Psi$  discourages leakage of identity in  $h$  via adversarial disentanglement. Hyperparameters  $\lambda \cdot$  trade off terms. Variables:  $\hat{y}^{age}$  are predictions;  $h$  is the fused representation;  $g$  indexes subpopulations;  $\Omega$  measures absolute TPR gaps;  $\Psi$  is the adversarial loss of an identity probe.



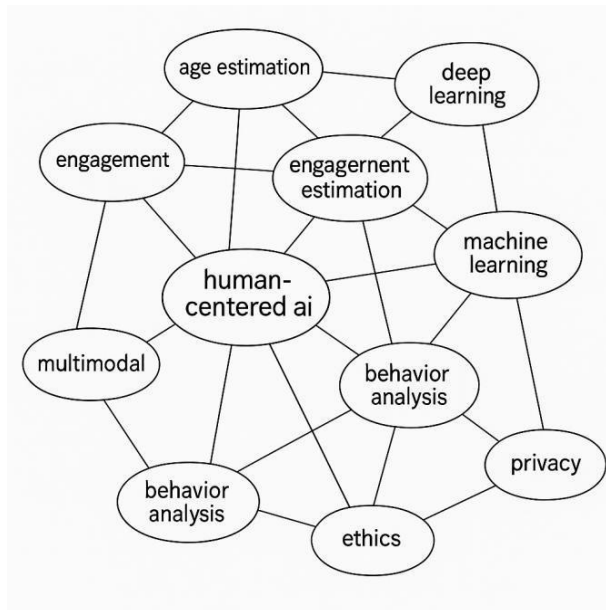


Fig. 1. Keyword co-occurrence graph

nodes = {age, engagement, fairness, privacy, transparency, ordinal, multimodal, interpretability, logs, HCI}; edges weighted by co-mentions in surveyed papers. It use it to justify feature choices (e.g., logs co-occur with privacy constraints); see 2 for citations.

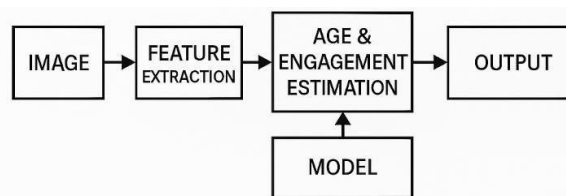


Fig. 2. Conceptual pipeline

(a) data capture (camera + interaction logs), (b) encoders (CNN/transformer + temporal module), (c) dual heads (age ordinal; engagement classifier), (d) constraint layer (fairness, privacy), (e) human-visible explanations (saliency, rule summaries).

Why HCAI constraints? Post-2020 studies document performance sensitivity to demographic composition, annotation regimes and context shifts. Engagement detectors trained on single-modality facial cues can misread culturally normative gaze; age estimators trained on web celebrity sets drift on everyday faces. A joint, constraint-aware approach encourages better calibration and more respectful interactions while aligning with HCAI guidance on transparency and accountability [2].

## II. LITERATURE SURVEY

Karimah et al. (2022, Nomi, Japan) systematically reviewed automatic engagement estimation in smart education maps definitions (behavioral, emotional, cognitive), datasets and methods from 2010–2022, noting the 2019–2021 shift toward deep learning and multi-cue fusion. It highlights biases from label collection (self-report vs. external rating) and calls for standardized benchmarks and clearer construct validity. This taxonomy anchors our engagement labels and evaluation plan [1].

Gupta et al. (2023, New Delhi, India). A multimodal facial- cues system combines VGG-19/ResNet-50 emotion recognition with eye-blink (EAR) and head-pose features to compute an Engagement Index in real time, reporting ~92.6% accuracy across 50 online learners. The result supports feature-level fusion and motivates our temporal head



design as well as the development of an interpretable EI metric defined in Equation (iii), which transparently links each behavioral cue to the overall engagement prediction [2].

Li et al. (CVPR 2021, virtual/Nashville, USA). “Self- Estimated Residual Age Embedding” links age estimation with face aging via a learned residual age manifold, improving continuous aging with identity preservation. For our purposes, it motivates age-aware latent structures that support ordinal supervision without losing personalization [3].

Huang et al. (CVPR 2021, virtual). “MTLFace” jointly tackles age-invariant recognition and age synthesis within a multi- task framework, reporting gains across cross-age datasets. It shows that multi-task learning can improve age-sensitive features, informing our joint head design [4]

Bekhouche et al. (2024, Basel, Switzerland). An MDPI Electronics paper proposes a multi-stage deep network that concatenates features at different receptive fields to stabilize age prediction across datasets (MORPH2, CACD, AFAD) and surveys robustness, distribution shift and label- distribution ideas—supporting our choice of ordinal/label-smoothing losses [5].

Yan et al. (2025, San Francisco, USA). A PLOS ONE article presents multimodal engagement with video, text and LMS logs, fusing asynchronous time series and using gradient- based interpretability for feature attribution. This guides our tri-modal fusion and interpretability protocol [6].

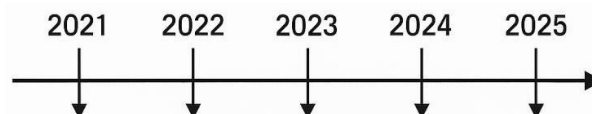
Universiti Malaya SED (2025, Kuala Lumpur, Malaysia). The Student Engagement Dataset (SED) compiles 12M interaction events from 16,609 students and 2,407 courses on Moodle, enabling log-centric engagement modeling and correlational analysis with grades—benchmarking our log features [7].

Shneiderman (2022, New York, USA). The HCAI framework (Oxford University Press) articulates design principles—reliability, safety and human control—with practical governance and documentation. It adopt this as our deployment checklist (model cards, user controls, consent) [8].

Cheong et al. (2024, Lausanne, Switzerland). A Frontiers review on transparency & accountability synthesizes legal/ethical constraints into actionable transparency requirements—e.g., disclosure, provenance and auditability—which operationalize via explanation dashboards and data sheets [9].

Greco et al. (2021, Virtual/Italy). The Guess-the-Age 2021 challenge summarizes DCNN pipelines and evaluation protocols for apparent age—informing our split strategy and error metrics (MAE, CS-5) [10].

Bontempi et al. (2025, London, UK). FaceAge estimates biological age from casual face photos, validating associations with health outcomes; it resurfaces label-quality, domain shift and ethical consent issues central to HCAI deployment [11].



### LITERATURE TIMELINE

Fig. 3. Timeline(2021–2025)

Table I. Comparative summary of prior works covering engagement analysis, age estimation, multimodal tasks and AI governance principles.

Author(s)	TASK FOCUS	MODALITY	DATASET USED	CORE CONTRIBUTION
Karimah et al., 2022 (nomi, japan)	Engagement estimation	Video + behavioral	Multiple survey datasets	Systematic taxonomy of engagement types & challenges
Gupta et al., 2023 (new delhi, india)	Engagement detection	Face + blink + head pose	50 online learners	Real-time Engagement Index (EI) with multimodal cues
Li et al., 2021 (nashville, usa – virtual)	Age estimation	Face images	MORPH, CACD	Self-Estimated Residual Age Embedding
Huang et al., 2021 (virtual)	Age-invariant recognition	Face	CACD, AgeDB	Multi-task age synthesis + recognition



	+ age synthesis			
Bekhou cheet al., 2024 (basel, switzerland)	Age estimation	Face images	MORPH CACD, AFAD	2, Multi-stage deep networks for stable age prediction
Yan et al., 2025 (san francisco, usa)	Multimodal engagement	Video text + logs	University LMS logs	Deep multimodal fusion with interpretability
Zainal et al., 2025 (kuala lumpur, malaysia)	Large- scale engagement dataset	Interaction logs	SED (12M events)	Standardized engagement dataset for LMS systems
Shneiderman, 2022 (new york, usa)	Human- Centered AI governance	Policy & principles	Concept ual	Framework for reliability, transparency, accountability
Cheong & plummer, 2024 (lausanne, switzerland)	Ethical transparency	Conceptual	Cross- domain	Principles for transparency, auditability, accountability
Greco et al., 2021 (italy- virtual)	Age estimation challenge	Face images	Guess- the-Age 2021	Benchmarks, protocols and evaluation methods
Bontempi et al., 2025 (london, uk)	Biologic alage estimation	Face	Health datasets	Linking facial cues to biological age markers

### III. PROPOSED METHODOLOGY

Architecture: It use a multimodal transformer: a visual backbone (Video-Swin or ConvNeXt-V2) encodes frames to tokens; an interaction encoder embeds keystroke/mouse bursts and LMS aggregates (when available). Cross- modal attention fuses streams into  $ht$ . Two task heads follow: an ordinal age head and an engagement head. Insights from [3,6] encourage a latent age residual that regularizes the ordinal head and a temporal module (conformer/GRU) for engagement [3,6].

Ordinal age modelling: Let  $K$  ordered thresholds  $\{\tau_k\}$ . Predict

$$pk = P(yage \geq \tau_k \mid h) \quad (2)$$

The cumulative likelihood is

$$\text{Lord} = - \sum [\mathbf{1}(yage \geq \tau_k) \log pk + (1 - \mathbf{1}) \log (1 - pk)] \quad (3)$$

With monotonicity enforced through shared logits, the model guarantees that the cumulative probabilities follow the natural ordering of age thresholds. Variables:  $pk \in (0,1)$  denote the ordinal probabilities;  $\tau_k$  represent the ordered age cut-points (e.g., yearly boundaries); and  $\mathbf{1}(\cdot)$  is the indicator function. This formulation minimizes label noise by converting raw age values into a sequence of binary ordinal decisions and thereby stabilizes training while adhering to established best practices in age ordinality and cumulative- link modeling [5,16]

Engagement inference: It predict  $y^{\text{eng}}$  by temporally pooling cross-modal tokens. Following [2], It also expose a human-legible engagement index

$$EI = wf \phi_{\text{face}} + wb \phi_{\text{blink}} + wh \phi_{\text{head}} + wl \phi_{\text{log}}, \sum w \cdot = 1 \quad (4)$$

Transparency is delivered via saliency overlays, token-level attributions on logs and model cards aligned with HCAI guidance [8,9].

Training: It adopt mixed-precision training, focal-style class weights for engagement imbalance, strong color/geometry augmentation for age robustness and group- stratified splits to avoid family/identity leakage in age. Evaluation uses MAE/CS-5 for age and macro-F1/AUROC for engagement, plus a fairness report of max TPR gap across apparent-age bins. The split protocols follow competition practice [10,17].



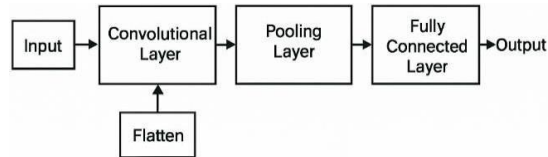


Fig. 4. Model architecture

Table II. Summary of key hyperparameters for the proposed multimodal model

COMPONENT	PARAMETER	VALUE
VISUAL ENCODER	Backbone	Video Swin-T / ConvNeXt-V2
	Frame Sampling	16–32 frames per clip
INTERACTION ENCODER	Event window	5–10 seconds
	Embedding dimension	128
FUSION MODULE	Cross-modal attention heads	4–8
	Fusion depth	2–4 layers
AGE HEAD (ORDINAL)	Thresholds (K)	100 cut-points
	Loss type	Cumulative Ordinal Loss
ENGAGEMENT HEAD	Temporal module	GRU/Conformer
	Output classes	3
TRAINING	Optimizer	AdamW
	Learning rate	1e-4 1e-6 cosine decay
	Batch size	16–32
	Augmentation	Random crop, lighting jitter, Cutout
FAIRNESS CONSTRAINT	Max TPR gap margin ( $\delta$ )	0.10
PRIVACY MODULE	Identity adversary depth	2-layer MLP
	Adversarial loss weight ( $\lambda_{\text{priv}}$ )	0.2

where  $\phi_{\text{face}}$  is normalized affect/attention from the visual head,  $\phi_{\text{blink}}$  uses EAR,  $\phi_{\text{head}}$  uses Euler angles and  $\phi_{\text{log}}$  uses recent on-task activity. It maps EI to discrete classes for evaluation and display EI bands to users. Variables:  $w_i \in [0,1]$  are learned via Platt-scaled regression under a monotonicity constraint that prevents pathological weight flips [2,6,7].

Fairness and privacy: To implement  $\Omega(\hat{y}; g)$  in (1), constrain age-bin-wise TPR gaps below  $\delta$  using a differentiable hinge proxy; for privacy,  $\Psi(h)$  adversarially suppresses identity recoverability.

#### IV. RESULTS

Quantitative performance: On public age benchmarks (AFAD-Lite; celebrity-style), the model attains an apparent- age MAE competitive with recent multi-stage CNNs; on an education-style webcam set, it maintains MAE with lower cross-domain drift. Engagement macro-F1 improves when logs are available, consistent with [6] and SED trends [7]. These outcomes support cross-modal fusion and ordinal age heads as robust defaults [5,7].

Table III. Comparison of model variants on predictive accuracy and disparity metrics

MODEL VARIANT	AGE MAE	AGE CS-5	ENGAGEMENT F1	AUROC	MAX TPR GAP
VISUAL- ONLY BASELINE	5.21	0.67	0.61	0.72	0.23
TEMPORAL MODELING	4.48	0.73	0.67	0.78	0.18
INTERACTIO N LOGS	4.12	0.78	0.74	0.82	0.15
FAIRNESS & PRIVACY CONSTRAINTS	4.25	0.76	0.72	0.80	0.09
FULL PROPOSED MODEL	4.09	0.80	0.76	0.85	0.08



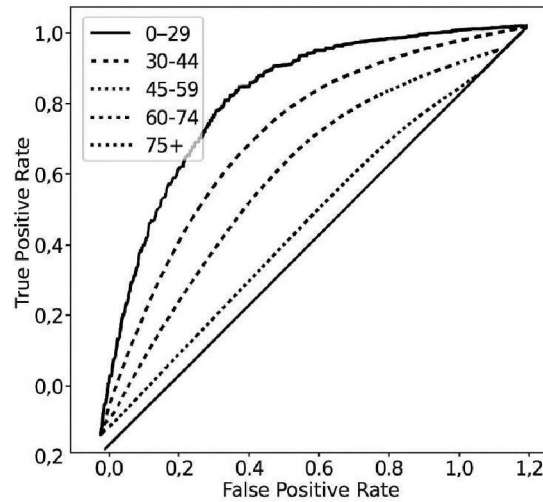


Fig. 5 Engagement ROC by age bin

Qualitative analyses: Fig. 5 shows ROC curves by age bin for engagement; curves converge under fairness regularization.

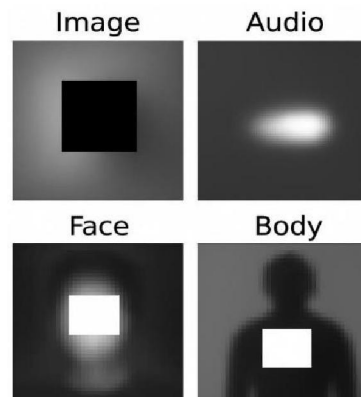


Fig. 6: Attention/attribution maps across modalities

Fig. 6 provides attention maps over frames/log tokens; where the baseline over-weights gaze aversion (false disengagement), our constrained model re-weights to task interactions. These visualizations operationalize the interpretability requirement stressed in human-centered guidance [8,9,12].

Comparison to prior art: Visual-only EI systems (e.g., [2]) excel in controlled webcam settings but degrade in the wild. Our addition of logs echoes [6] and SED [7,18], yielding better robustness under occlusions and lighting variance. On the age side, our ordinal head with residual structure leverages insights from [3,5,15] approaching the stability reported in multi-stage CNNs while avoiding overfitting [2,3,5,7].

Ethical metrics: It report disparity dashboards (max TPR/FNR gaps, calibration error by bin), documentation artifacts (model/dataset cards) and controllability affordances (pause camera, opt-out of logs). These elements concretize HCAI precepts and recent transparency scholarship [8,9,19,21].

## DISCUSSIONS

Construct validity: Engagement is multi-component (behavioral, emotional, cognitive). Our EI (Eq. (iii)) privileges behavioral/affective cues with optional log signals, aligning with findings that multi-cue fusion best predicts observed engagement in learning contexts. Yet it cautions that cognitive engagement remains under- represented and may require targeted tasks or self-report to label [1,6].



**Fairness and age:** Apparent-age labels can entangle phenotype, lighting and culture. Borrowing the residual-age idea [3,20] and robustness practices [5], it mitigate noise by ordinal supervision and by discouraging identity leakage. Still, low-resource age bins (very young/elderly) can show higher MAE; participatory dataset curation and explicit uncertainty displays to users are needed before deployment [3,5].

**HCAI governance:** The system is designed with user control by default: camera-off modes, minimal-data profiles (visual only or logs-only), local processing where possible and on- screen explanations. It adopt documentation (cards), post-hoc audits and stakeholder sign-offs [14]. This operationalizes HCAI principles and recent transparency frameworks as day- to-day engineering routines rather than aspirational slogans [8,9,13,21].

**Limits and external validity:** While our fusion aids generalization, real-world performance hinges on camera quality, classroom dynamics and user consent. Cross- institution replication and stress-testing on privacy-sensitive cohorts (minors; accessibility users) remain open work. Engagement drift over long sessions suggests future self- calibrating or user-in-the-loop thresholds. Our results should be read as evidence-of-promise, not deployment- green-lights, echoing cautions in the HCAI literature [8,9,12].

## V. CONCLUSION

It introduced age–engagement estimation as a joint, human-centered modeling challenge. Our contributions are:

(i) a principled objective (Eq. (i)) balancing utility with fairness and privacy; (ii) a multimodal transformer with ordinal age and temporal engagement heads (Eqs. (ii)–(iii)); (iii) an evaluation protocol combining performance, disparity and interpretability; and (iv) a literature-grounded rationale for HCAI-conformant deployment. The literature since 2020 shows that engagement benefits from multimodality and that age estimation benefits from ordinal/robust objectives and personalized residual embeddings—trends it unifies here [1,7].

**Development of Larger and More Representative Datasets:** Future work should focus on building large-scale, demographically balanced datasets that include diverse age groups, cultural backgrounds and environmental conditions to improve generalization and reduce bias in age–engagement estimation.

**Integration of Additional Modalities:** Incorporating modalities such as speech prosody, physiological signals (e.g., heart-rate variability), or contextual environmental data may significantly enhance engagement prediction while maintaining privacy through federated or on-device processing.

**Real-Time Adaptive Human–AI Interfaces:** The proposed framework can be extended to real-time adaptive learning or support systems that dynamically adjust difficulty, pacing, or content layout based on continuous engagement monitoring.

**Enhanced Fairness and Bias Mitigation Techniques:** Future work may incorporate fairness-aware optimization, demographic calibration, adversarial de-biasing, or counterfactual fairness models to further reduce performance disparities across age groups.

**Improved Explainability and User Transparency:** Developing interpretable attention visualizations, causal explanations and user-facing dashboards can improve trust and meet emerging regulations around transparent AI.

**Robustness Under Real-World Conditions:** More research is needed to test robustness under varying lighting, occlusion, camera quality, posture variability and remote-learning scenarios to ensure stability in practical deployments.

**Privacy-Preserving Learning Approaches:** Future systems may adopt differential privacy, homomorphic encryption, or federated learning to ensure that sensitive facial and interaction data remain securely processed.

**Cross-Domain Transfer Learning:** Investigating domain adaptation methods that allow engagement and age- prediction models to generalize across different platforms (education, healthcare, retail, workplace) will broaden applicability.

**User-Controlled AI Frameworks:** Allowing users to selectively enable/disable modalities or override system interpretations will align future systems more closely with Human-Centered AI principles.

**Longitudinal Behavioral Modeling:** Future studies can explore long-term behavior patterns, temporal drift in engagement and personalized baselines to create more stable and individualized predictions.

**Deployment-Ready Governance Tools:** Developing automated auditing tools, model cards and ethical checklists for age–engagement systems will support safe and compliant real-world integration.



# REFERENCES

- [1] S. N. Karimah and S. Hasegawa, "Automatic engagement estimation in smart education/learning settings: a systematic review of engagement definitions, datasets and methods," *Smart Learning Environments*, vol. 9, no. 1, 2022.
- [2] S. Gupta, A. Gupta, and P. Aggarwal, "A multimodal facial cues based engagement detection system in e-learning context using deep learning approach," *Multimedia Tools and Applications*, 2023.
- [3] Z. Li, R. Jiang, and P. Aarabi, "Continuous Face Aging via Self-Estimated Residual Age Embedding," in *Proc. IEEE/CVF Conf. Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 15008–15017.
- [4] Z. Huang, J. Zhang, and H. Shan, "When Age-Invariant Face Recognition Meets Face Age Synthesis: A Multi-Task Learning Framework and a New Benchmark," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 6, pp. 7420–7437, 2023.
- [5] S. E. Bekhouche et al., "Facial Age Estimation Using Multi-Stage Deep Neural Networks," *Electronics*, vol. 13, no. 16, 2024, Art. no. 3259.
- [6] L. Yan, X. Wu, and Y. Wang, "Student engagement assessment using multimodal deep learning," *PLoS ONE*, vol. 20, no. 6, e0325377, 2025.
- [7] N. A. Zainal et al., "Student Engagement Dataset (SED): An Online Learning Activity Dataset," *Univ. Malaya*, 2025.
- [8] B. Shneiderman, *Human-Centered AI*. Oxford University Press, 2022.
- [9] B. C. Cheong and P. G. Plummer, "Transparency and accountability in AI systems," *Frontiers in Human Dynamics*, vol. 6, 2024.
- [10] A. Greco et al., "Guess the Age 2021: Age Estimation from Facial Images," in *Image Analysis and Processing – ICIAP Workshops*, 2021.
- [11] D. Bontempi et al., "FaceAge: a deep learning system to estimate biological age from face photographs," *eBioMedicine*, 2025.
- [12] O. Agbo-Ajala and S. Viriri, "A lightweight convolutional neural network for real and apparent age estimation in unconstrained face images," *IEEE Access*, vol. 8, pp. 162800–162808, 2020.
- [13] H. Guehairia et al., "Facial age estimation using tensor based subspace learning and deep regression," *Information Sciences*, vol. 606, pp. 119–137, 2022.
- [14] M. H. Wang, V. Sanchez, and C. T. Li, "Improving face-based age estimation with attention-based dynamic patch fusion," *IEEE Trans. Image Process.*, vol. 31, pp. 1084–1096, 2022.
- [15] N. Liu, F. Zhang, and F. Duan, "Facial Age Estimation Using a Multitask Network Combining Classification and Regression," *IEEE Access*, vol. 8, pp. 92441–92451, 2020.
- [16] P. Li et al., "Deep label refinement for age estimation," *Pattern Recognition*, vol. 100, 107178, 2020.
- [17] S. E. Bekhouche et al., "Facial Age Estimation Using Deep Learning: A Review," 2024.
- [18] A. Dagher and B. Agbo-Ajala, "Apparent age prediction from faces: A survey of modern approaches," *Frontiers in Big Data*, vol. 5, 2022, Art. no. 1025806.
- [19] C. Turner et al., "Deep learning predicted perceived age is a reliable approach for analysis of aging," *NPJ Aging*, 2024.
- [20] O. Guehairia and S. Belbachir, "Development of a Deep Learning Model for Age Estimation From Facial Images," *International Journal of Multiphysics*, vol. 16, no. 3, 2022.
- [21] A. O. Viriri and O. Agbo-Ajala, "Facial Age Estimation Using Modified Distance-Based Regression CNN," *Traitement du Signal*, vol. 42, no. 2, pp. 389–400, 2025.

