

# Multimodal Machine Learning in Mental Health: A Survey of Data, Algorithms, and Challenges

ZAHRAA AL SAHILI, Queen Mary University of London, United Kingdom

IOANNIS PATRAS, Queen Mary University of London, United Kingdom

MATTHEW PURVER, Queen Mary University of London & Jožef Stefan Institute, United Kingdom & Slovenia

Multimodal machine learning (MML) is rapidly reshaping the way mental-health disorders are detected, characterized, and longitudinally monitored. Whereas early studies relied on isolated data streams—such as speech, text, or wearable signals—recent research has converged on architectures that integrate heterogeneous modalities to capture the rich, complex signatures of psychiatric conditions. This survey provides the first comprehensive, clinically grounded synthesis of MML for mental health. We (i) catalog 26 public datasets spanning audio, visual, physiological signals, and text modalities; (ii) systematically compare transformer, graph, and hybrid based fusion strategies across 28 models, highlighting trends in representation learning and cross-modal alignment. Beyond summarizing current capabilities, we interrogate open challenges—data governance and privacy, demographic and intersectional fairness, evaluation explainability, and mental health disorders complexity in multimodal settings. By bridging methodological innovation with psychiatric utility, this survey aims to orient both ML researchers and mental-health practitioners toward the next generation of trustworthy, multimodal decision-support systems.

**CCS Concepts:** • Computing methodologies → Machine learning; • Applied computing → Life and medical sciences.

**Additional Key Words and Phrases:** Multimodal machine learning, mental health, healthcare

## ACM Reference Format:

Zahraa Al Sahili, Ioannis Patras, and Matthew Purver. 2025. Multimodal Machine Learning in Mental Health: A Survey of Data, Algorithms, and Challenges. 1, 1 (June 2025), 12 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

## 1 Introduction

Mental disorders constitute a larger share of global disability than any other non-communicable disease group. In 2019 roughly 970 million people—one in eight worldwide—were living with clinically significant anxiety, depression, or a related condition [1]. Treatment capacity has not kept pace: up to 85 % of affected individuals in low- and middle-income countries receive no formal care, and even in high-income settings the median delay from symptom onset to first intervention exceeds a decade [2]. Bridging this chasm demands solutions that are simultaneously scalable, affordable, and clinically trustworthy.

Digital platforms are beginning to fill that void. The global market for web- and app-based mental-health services was valued at \$5.1 billion in 2020 and is projected to grow at more than 20 % annually through 2028 [3]. Beyond tele-therapy and self-guided interventions, *multimodal* machine learning (MML)—the computational fusion of heterogeneous signals such as language, prosody, facial micro-expressions, kinematics, and physiology—has

---

Authors' Contact Information: Zahraa Al Sahili, Queen Mary University of London, United Kingdom, z.alsahili@qmul.ac.uk; Ioannis Patras, Queen Mary University of London, United Kingdom, i.patras@qmul.ac.uk; Matthew Purver, Queen Mary University of London & Jožef Stefan Institute, United Kingdom & Slovenia, m.purver@qmul.ac.uk.

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM XXXX-XXXX/2025/6-ART

<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

emerged as an especially promising avenue. A growing body of evidence shows that combining complementary cues yields markedly higher diagnostic and prognostic accuracy than any single channel alone; for example, jointly modelling vocal inflection, facial affect, and heart-rate variability can uncover depressive symptomatology that text-only or audio-only models routinely miss [4].

Yet translating this promise into routine clinical workflows is non-trivial. Fragmented data silos, inconsistent annotation schemes, algorithmic bias, privacy constraints, and a scarcity of longitudinal, ecologically valid benchmarks all impede deployment. An up-to-date, critical synthesis of the literature is therefore both timely and necessary.

*Scope and contributions.* This survey organises and critiques the rapidly expanding work on multimodal ML for mental health:

- (1) **Data landscape** – We catalogue 26 public corpora, ranging from laboratory interviews to in-the-wild smartphone and wearable streams, and map each dataset to clinical targets such as major depressive disorder, post-traumatic stress disorder, stress, and bipolar disorder.
- (2) **Modelling approaches** – We synthesise the state of the art across representation learning, fusion strategies (hybrid, graph and attention).
- (3) **Open challenges** – We analyse persistent obstacles—data, privacy, fairness, evaluation explainability, and mental health disorders complexity.

By clarifying both the promise and the pitfalls of multimodal ML, we aim to orient machine-learning researchers, clinicians, and policymakers toward solutions that are methodologically sound, ethically grounded, and clinically actionable, ultimately advancing the scale and quality of global mental-health care.

## 2 Data Types and Datasets

Multimodal learning for mental-health research rests on two pillars:

- (1) complementary data streams that capture behaviour and physiology at different levels of abstraction, and
- (2) curated corpora that allow models to be trained, validated and fairly compared.

This section surveys both, providing a taxonomy of signal types (§2.1) and an inventory of benchmark datasets (§2.2).

### 2.1 Data modalities

Text captures cognitive content and latent affect; audio supplies paralinguistic prosody that often precedes explicit reports; video provides non-verbal behaviour such as micro-expressions or psychomotor retardation; and physiological channels offer objective indices of autonomic or central-nervous-system activity. Fusing these streams lets models disentangle overlapping symptom profiles and improves robustness to missing or noisy channels.

### 2.2 Benchmark multimodal mental-health datasets

Table 2 summarises the 24 public corpora most frequently cited in the literature ([5]–[30]). Together they span six diagnostic themes—depression, stress, PTSD, bipolar disorder, behavioural disorders and generic emotion/affect recognition—and four primary modalities.

*Diagnostic coverage.* Depression and stress dominate (7 and 10 datasets, respectively), reflecting their global prevalence and the relative ease of eliciting symptoms in laboratory tasks. PTSD and bipolar disorder remain under-represented, while schizophrenia and other psychotic disorders are largely absent.

Table 1. Typical sources, salient features and common clinical tasks for each modality.

| Modality          | Typical sources  | Salient features  | Common tasks  |
|-------------------|--|---|---|
| <b>Text</b>       | social-media posts, EHR notes, therapy transcripts, smartphone messages  | lexical affect, syntactic complexity, semantic coherence, pronoun use         | screening; relapse prediction; suicidal-ideation detection                    |
| <b>Audio</b>      | structured interviews, telephone calls, vlog diaries, ambient sound      | prosody (F0, intensity), speech rate, voice quality, spectral entropy         | depression screening; bipolar-episode detection; stress recognition           |
| <b>Video</b>      | webcam or smartphone recordings, clinical assessments, in-the-wild vlogs | facial action units, gaze, head pose, body kinematics                         | emotion recognition; PTSD detection; severity rating                          |
| <b>Physiology</b> | wearables (PPG, EDA, accelerometers), EEG/ERP, fMRI, eye-tracking        | heart-rate variability, skin conductance, EEG power bands, pupillary response | affect recognition; anxiety & stress quantification; sleep-quality assessment |

*Modality combinations.* Six corpora include the full triad of video–audio–text; another six provide paired video–audio recordings. Physiology appears in 13 datasets but only rarely alongside language, leaving room for richer, ecologically valid sensor fusions.

*Scale and demographics.* Sample sizes range from 25–275 laboratory participants to tens of thousands of utterances in movie/dialogue corpora. Only nine datasets report detailed gender balance and even fewer disclose ethnicity or socio-economic status, hampering fairness analysis.

*Access and ethics.* All corpora are gated by data-use agreements; several require proof of IRB clearance. Researchers should budget time for approvals and verify that planned tasks align with original consent scopes.

### 2.3 Common patterns and persistent gaps

Only one-quarter of datasets combine *all* behavioural channels with any physiology, limiting research on end-to-end biopsychosocial models. Label quality ranges from validated scales (PHQ-8/9, HAMD, MADRS) in clinical interviews to self-report hashtags in social-media corpora, complicating cross-dataset evaluation. Longitudinal depth is rare—most datasets provide a single session per participant—and geographic skew favours Europe and North America. Broader cultural representation and bias-aware splits are therefore essential for robust generalisation.

### 2.4 Implications for future collection

Progress now hinges on resources that mirror real life: culturally diverse cohorts, device-agnostic capture, repeated measures over months rather than hours and labels grounded in clinician assessment *and* self-report. Promising directions include federated frameworks where raw data never leave the collection site and integrated digital phenotyping that merges smartphone behaviour, ecological momentary assessment and wearable physiology.

Table 2. Publicly available multimodal datasets for mental-health research. Modalities: **V**=Video, **A**=Audio, **T**=Text, **P**=Physiology. “–” means not reported.

| Ref. | Dataset                        | Target disorder(s) | Modalities <sup>†</sup> | N<br>subj./sessions  | Demographics | Country /<br>Source   |
|------|--------------------------------|--------------------|-------------------------|----------------------|--------------|-----------------------|
| [5]  | DAIC-WOZ                       | Depression, PTSD   | V A T                   | 189 sessions         | –            | USA (lab)             |
| [6]  | E-DAIC                         | Depression, PTSD   | V A T                   | 275 subj. / 70 h     | –            | USA (lab)             |
| [7]  | AVEC 2013                      | Depression         | V A                     | 340 videos           | –            | Europe<br>(challenge) |
| [8]  | Multi-Modal<br>Mental-Disorder | Depression         | A T P                   | 52–55<br>subj./state | –            | China                 |
| [9]  | SWELL                          | Stress             | P                       | 25 subj.             | 32% F        | Netherlands           |
| [10] | D-Vlog                         | Depression         | V A                     | YouTube clips        | –            | YouTube               |
| [11] | CMDC                           | Depression         | V A                     | 45 subj.             | –            | China                 |
| [12] | BPC+SEWA+RECOLA                | Bipolar, behaviour | V A T P                 | merged               | –            | TR/DE/HU/FR           |
| [13] | Turkish BDC                    | Bipolar            | V A                     | 51 subj.             | 31% F        | Turkey                |
| [14] | PTSD-in-the-Wild               | PTSD               | V A T                   | 634 subj.            | –            | USA                   |
| [15] | BIRAFFE2                       | Emotion            | T P                     | 103 subj.            | 33% F        | Poland                |
| [16] | DEAP                           | Emotion            | V P                     | 32 subj.             | –            | UK                    |
| [17] | MuSE                           | Stress             | V T P                   | 28 subj.             | 32% F        | USA                   |
| [18] | MIREX                          | Emotion            | V A T                   | 193 samples          | –            | –                     |
| [19] | MELD                           | Emotion (dialogue) | V A T                   | 13k utterances       | –            | Movie subtitles       |
| [20] | WEMAC                          | Emotion            | V A T P                 | 47 women             | 100% F       | Spain                 |
| [21] | EmoReact                       | Emotion (children) | V                       | 63 children          | 51% F        | YouTube               |
| [22] | WESAD                          | Stress & affect    | P                       | 25 subj.             | –            | Switzerland           |
| [23] | Nurse Stress                   | Stress             | P                       | 15 nurses            | –            | USA                   |
| [24] | MuSe-Stress                    | Stress             | V A                     | 105 subj.            | 70% F        | Germany               |
| [25] | Multimodal Stress              | Stress             | V A T P                 | 80 subj.             | 59% F        | –                     |
| [26] | EmpathicSchool                 | Stress             | V P                     | 20 subj.             | –            | Finland & USA         |
| [27] | Self-Adaptors                  | Stress             | V A T                   | 35 subj.             | –            | USA                   |
| [28] | CLAS                           | Stress             | T P                     | 62 subj.             | –            | Singapore             |
| [29] | MuSe 2022                      | Stress             | V A                     | 105 subj.            | –            | UK/Germany            |
| [30] | UBFC-PHYS                      | Stress             | P                       | 56 subj.             | –            | France                |

<sup>†</sup> Modalities key: **V** = Video, **A** = Audio, **T** = Text, **P** = Physiology.

Such investments are costly and ethically complex, yet without them the gains achieved on current benchmarks risk stalling at the point of clinical deployment.

### 3 Methods

Deep learning has reshaped multimodal mental-health modelling. Where early work stitched together hand-engineered features with support-vector machines or random forests, contemporary studies rely on three neural families—convolutional-recurrent hybrids, transformers, and graph neural networks (GNNs)—each tuned to a different aspect of the fusion challenge. Convolutional-recurrent models thrive on compact, well-synchronised recordings; transformers scale to loosely aligned, irregular streams; and GNNs turn the tangled relationships

among cues into a source of predictive power. The remainder of this section traces how these families complement one another, starting with the still-indispensable convolutional–recurrent pipelines.

### 3.1 Hybrid CNN/RNN methods

Hybrid pipelines that pair *stream-specific* convolutional or recurrent encoders with lightweight cross-modal fusion dominated the first wave of multimodal mental-health work. Each modality is processed by a network tailored to its signal characteristics—CNNs for images and spectrograms, RNNs for word sequences or temporal features—after which the per-stream vectors are merged by concatenation, gating, or simple voting. Although later eclipsed by transformers and GNNs, these CNN/RNN systems still provide strong, computationally efficient baselines across a wide range of disorders.

Early clinical-interview research illustrates the value of *temporal* hybrids. In DAIC-WOZ, highway gates first suppress noisy audio-visual frames, then concatenate the filtered cues with GloVe text embeddings before an LSTM; the scheme pushes the depression-screening  $F_1$  to 0.81 and halves PHQ-8 error over naïve concatenation [31]. A richer three-branch design for AVEC fuses RGB faces with “motion-code” images and dual audio streams inside CNN–BiLSTM encoders, then applies hierarchical attention across modalities; the cascade trims MAE to 6.48–7.01 and beats pure early or late fusion on both AVEC benchmarks [32].

On Twitter timelines, signal-to-noise is the bottleneck. COMMA tackles this by training two reinforcement-learning agents—a text selector and an image selector—that filter posts before a GRU + VGG early fusion; Macro- $F_1$  climbs to 0.90 and remains robust even when depressed users are only 10 % of the pool [33]. A similar “filter-then-fuse” logic appears on Instagram: a fine-tuned BERT and a 10-layer CNN make independent predictions that are merged by 70:30 soft voting, delivering 99 % accuracy and  $F_1$  on a 10 k-post corpus [34].

Hybrid CNN/RNN stacks also reach less studied pathologies. For post-partum depression, AlexNet features from questionnaire text and Mel-spectrograms are early-concatenated and fed to an attentive Bi-LSTM, raising  $F_1$  above 0.98 on the UCI PPDD corpus [35]. Psychological stress in WESAD is handled by mapping four physiological streams to wavelet images, extracting channel-attentive SqueezeNet features, pruning them with an arithmetic optimiser and classifying via a DenseNet-LSTM; the full pipeline attains 99 % accuracy and 97.8 % Macro- $F_1$  [36].

Even single-sensor scenarios can profit from hybrid thinking. By viewing MFCCs and Spectro-CNN embeddings as “pseudo-modalities,” early concatenation plus a deep fully-connected classifier drives audio-only depression models to  $\approx 90$  % accuracy in English and Mandarin corpora [37]. Conversely, [38] shows that a single face still harbours multiple diagnostic hints: emotion probabilities from a YOLO detector are mapped to disorder labels, and penultimate embeddings from MDNet, ResNet-50 and a ViT are concatenated then soft-maxed, yielding 81 % accuracy while remaining explainable through Grad-CAM.

The corpus of hybrid CNN/RNN work yields several broad insights that remain relevant even as the field pivots toward transformers and graphs. First, *when* fusion occurs matters: early concatenation shines when signals are synchronous and clean, whereas late voting or RL-based selection excels under heavy noise or redundancy. Second, attention and gating—precursors of transformer cross-modal layers—consistently boost performance by spotlighting reliable channels and suppressing artefacts. Third, these CNN/RNN hybrids offer favourable speed–accuracy trade-offs, making them attractive for on-device inference or clinical settings with limited compute. Their main limitations are small dataset sizes, scarce fairness analysis and limited capacity to capture very long-range dependencies—gaps that have since motivated the rise of transformer and GNN architectures explored in the next subsections.

### 3.2 Transformer-based multimodal methods

Transformers have rapidly become the work-horse for multimodal mental-health research because the self-attention mechanism offers a uniform way to integrate heterogeneous cues while preserving long-range dependencies. Recent studies illustrate how this architecture family has been adapted—sometimes ingeniously—to fuse speech, language, vision, physiology and even genomic data for diagnosis or risk screening.

Early attempts centred on *in-the-wild* social-media content. A cross-attention transformer trained on the audio–visual *D-Vlog* corpus combines eGeMAPS speech descriptors with facial–landmark trajectories; its bi-directional attention layers let acoustic queries attend to visual keys and vice-versa, lifting the  $F_1$ -score to 63.5 % on a 961-video benchmark and beating BLSTM and tensor-fusion baselines by up to four points [10]. Moving from videos to timelines, a multimodal LXMERT-style encoder enriched with *time2vec* positional codes fuses CLIP image tokens and EmoBERTa sentence embeddings at the post level; the design secures an  $F_1$  of 0.931 on Twitter and 0.902 on Reddit—two to five points above earlier GRU + VGG systems while still scaling to millions of posts [39]. A companion study on extended vlogs pushes further, inserting video patches, wav2vec-2.0 audio frames and Whisper-BERT transcripts into a single 12-layer vision–language transformer; this early, latent-space fusion gains 4.3  $F_1$  points over the best cross-attention CNN and even transfers competitively to the clinical DAIC-WOZ set [40]. In text-only work, six domain-specific BERT variants detect depression and suicidality across four Reddit/Twitter corpora; although unimodal, they still exemplify transformer fusion because lexical, syntactic and discourse cues are integrated inside multi-head self-attention, reaching up to  $F_1 = 0.967$  [41].

Clinical interviews motivate tighter control of temporal structure. The *Multimodal Purification Fusion* network extracts EfficientNet-BiLSTM acoustic features and BiLSTM sentence embeddings, then decomposes each into modality-specific and modality-common parts before a cross-attention transformer recombines them; the method attains  $F_1 = 0.88$  on DAIC-WOZ, three to five points over naïve concatenation and audio- or text-only baselines [42]. A lighter *Topic-Attentive Transformer* keeps RoBERTa and wav2vec encoders separate until a late concatenation, but gates the textual stream with learned topic weights derived from ten canonical interview questions; despite its simplicity, the model still records  $F_1 = 0.647$  and shows that selectively emphasising clinically salient segments can compensate for small sample sizes [43].

Transformer variants also prove effective outside speech and language. *DepMSTAT* introduces a two-stage *Spatio-Temporal Attentional Transformer*: the first block attends across facial–landmark dimensions within each frame, the second along time, before a cross-modal attention fuses video with audio; precision–recall balances of roughly 73–76 % on the enlarged D-Vlog dataset underscore the benefit of separating spatial from temporal attention [44]. In physiological monitoring, *MUSER* concatenates BERT text embeddings and eGeMAPS audio features, then trains the shared transformer under an adaptive multi-task curriculum that balances stress labels with auxiliary arousal/valence signals; dynamic sampling yields  $F_1 = 86.4$  %, markedly above static curricula and late-fusion GRUs [45]. Finally, precision psychiatry enters the transformer arena through a ViT + XGBoost ensemble that sums logits from brain-MRI patches and polygenic-risk scores; the hybrid lifts area-under-the-curve to 0.891—a 0.216 gain over genetics alone and 0.068 over imaging alone—highlighting how self-attention can mine neuro-anatomical patterns even at modest sample sizes [46].

Taken together, these works reveal several insights. *Cross-attention* has emerged as the dominant fusion mechanism when modalities are aligned in time, whereas *early latent fusion* with homogeneous self-attention suits settings where synchrony is weak or absent. Pre-training—on ImageNet for ViT, on large speech corpora for wav2vec 2.0, and on domain-specific Reddit dumps for MentalBERT—systematically boosts downstream accuracy, mitigating the chronic scarcity of mental-health labels. Temporal encoding remains an open question: *time2vec* helps on dense timelines, but permutation-invariant SetTransformers win when postings are sporadic. Explainability is now expected, whether via SHAP word saliences [41] or Grad-CAM maps over facial regions [44]. Yet challenges persist: most datasets still contain fewer than one thousand subjects; cross-cultural robustness

and fairness across gender or dialect are rarely audited; and privacy constraints limit access to richer modalities such as EEG or smartphone sensors. Even so, the transformer toolkit—flexible fusion, scalable pre-training and growing interpretability—has already advanced state-of-the-art performance by 2–10 percentage points across domains, signalling its central role in the next generation of multimodal mental-health technology.

### 3.3 Graph-neural approaches

Graph neural networks (GNNs) have emerged as a natural fit for multimodal mental-health research because they can encode heterogeneous entities—brain regions, interview segments, social-media objects—as nodes and let information flow along richly typed edges. The most recent studies, spanning clinical interviews, social-media videos, resting-state fMRI and conversational corpora, show how graph reasoning can be moulded to fuse modalities, tame data scarcity and expose interpretable biomarkers.

In the clinical-interview setting, three works inject graph structure *after* a sequential encoder so that sparse subjects can borrow strength from one another. A knowledge-aware graph-attention network couples audio, video and text nodes with domain-defined meta-paths, then gates the result with a temporal-convolutional network, pushing DAIC-WOZ  $F_1$  to 0.95–28 points above classical early fusion [47]. A few-shot variant first learns modality weights through Bi-LSTM pre-fusion, then constructs a fully connected support-query graph whose edges are refined by message passing; this raises accuracy from 72 % (plain Bi-LSTM) to 86 % [48]. Extending to the larger E-DAIC corpus, a heterogeneous graph transformer that links intra- and inter-modal chunks attains 4.8 RMSE on PHQ-8 and still transfers with 78 %  $F_1$  to an external coaching dataset [49].

Graph thinking also broadens depression detection beyond the face-centred paradigm of earlier vlogging work. In MOGAM, every YouTube frame is parsed by YOLO into COCO objects; co-occurrence graphs are fused with ResNet scene embeddings and KoBERT metadata through transformer cross-attention, lifting  $F_1$  to 0.888 for daily-vs-depressed vlogs and hitting 0.997 for high-risk detection, while zero-shot transfer to the English D-Vlog set remains competitive [50].

Neuro-imaging studies push GNNs into the realm of precision psychiatry. Treating each brain as a functional-connectivity graph, an unsupervised graph auto-encoder plus FCNN reaches 72 % accuracy on a small Duke-MDD cohort and still tops competing CNNs on the 477-subject REST-meta-MDD repository [51]. MAMF-GCN builds a population graph with two parcellations and phenotype edges; channel-common convolution and attention fusion almost saturate the Southwestern MDD (99.2 % acc.) and ABIDE-ASD (97.7 % acc.) benchmarks [52]. At an even finer temporal scale, MTGCN learns graph-convolution attention inside multi-atlas dynamic FC windows, then merges them via a multimodal transformer, achieving 81 % accuracy for insomnia disorder while highlighting dysfunctional DMN regions [53]. A sister study on ABIDE I shows that three feature-scale population sub-graphs processed by deep Chebyshev GCNs and concatenated embeddings can push ASD screening to 91.6 % accuracy and 95.7 % AUC, a 12–20 pp jump over earlier graph pipelines [54].

Conversation-level emotion recognition—viewed as an upstream proxy for mood monitoring—also benefits from GNN fusion. MMGCN places audio, video and text nodes from each utterance in one graph so that spectral convolutions interleave modalities; it delivers 66.2 %  $F_1$  on IEMOCAP and 58.7 % on MELD, edging recurrent and transformer competitors [55]. COGMEN further augments this with relational edges for speaker turns and temporal direction; its GraphTransformer lifts IEMOCAP-6  $F_1$  to 67.6 % and MELD to 58.7 %, with ablations confirming the indispensability of distinct edge types [56].

Language-only risk screening shows that GNNs can add value even when no explicit second modality exists. MM-EMOG casts the entire Twitter or Reddit corpus as a word-document graph whose edges encode multi-label emotion co-occurrence; a two-layer GCN fine-tuned on BERT boosts  $F_1$  by 8–21 points over transformer baselines across suicide and depression sets while retaining full privacy by ignoring user histories [57].

Across these studies several patterns surface. First, *where* fusion happens varies: some models mix modalities inside each graph convolution [50, 55], others concatenate graph embeddings with CNN or transformer features only at the final layer [51]. Second, heterogeneity is key: typed edges (speaker roles, phenotypes, object co-occurrences) consistently raise accuracy by 2–6 pp over homogeneous graphs. Third, population graphs unlock learning in low-sample regimes, but require careful edge design—phenotype-weighted similarity [54] or GNN-learned adjacency [52] both outperform naïve  $k$ -NN. Finally, several works expose interpretable biomarkers: DMN-centred attention maps for insomnia [53], middle-occipital connectivity for MDD [51], and edge saliences for suicide lexicons [57], underscoring GNNs’ potential to bridge black-box performance with clinical insight.

## 4 Open Challenges

Despite a surge of methodological innovation, translating multimodal machine-learning research into routine mental-health practice remains fraught with obstacles. These span the entire development pipeline—from data acquisition to model governance—and are intertwined in ways that call for multidisciplinary solutions rather than purely technical fixes.

### 4.1 Data availability

Multimodal models prosper only when they can observe the full biopsychosocial spectrum of human experience, yet the field still relies on a handful of convenience samples drawn largely from Western, university-educated populations. Privacy regulations such as GDPR and HIPAA, while vital, make it costly to gather and share raw video, audio and physiology; marginalised communities are often excluded because obtaining IRB approval, culturally appropriate consent and secure storage is harder in resource-constrained settings. The result is a patchwork of small, siloed datasets that encourage over-fitting and complicate external validation. Federated learning, synthetic data and privacy-preserving cryptography are promising but immature countermeasures; they require robust proof that utility is not sacrificed at the altar of privacy.

### 4.2 Benchmarks and reproducibility

Meaningful progress depends on comparing algorithms under identical conditions, yet today’s benchmark landscape is fragmented. Studies vary in how they split data, which severity scales they predict and whether they evaluate at the clip, session or patient level. A reproducible suite should mirror real-world deployment: it must include held-out sites, force models to contend with missing or corrupted modalities and report clinically salient metrics such as calibration, positive predictive value at low prevalence and time-to-detection. Without such guard-rails, incremental gains on familiar leaderboards tell us little about clinical utility.

### 4.3 Explainability and clinical trust

Mental-health assessment is as much an interpretive craft as it is a measurement science. Clinicians are unlikely to act on black-box scores that cannot be reconciled with observable signs or patients’ self-reports. Yet the very factors that make multimodal models powerful—high-dimensional fusion, cross-attention, recursive graph-reasoning—tend to obscure causal pathways. Post-hoc methods such as SHAP, LIME or Grad-CAM offer local saliency but seldom convey longitudinal or cross-patient consistency, and they can mislead when features are correlated. More promising are intrinsically interpretable architectures that align attention with DSM-5 symptom clusters, propagate evidence along clinically meaningful graphs or generate counterfactual narratives that clinicians can critique.

#### 4.4 Bias and fairness

Because training corpora under-represent certain age groups, genders, ethnicities and dialects, multimodal models risk perpetuating the very disparities they purport to alleviate. Bias can creep in through sensor placement (dark-skinned faces are harder to track), language variation (AAVE mislabelling) or differential access to care (labels mirror systemic inequities). Mitigation therefore demands a full-lifecycle approach: equitable data collection, algorithmic debiasing (e.g. adversarial or counterfactual training) and post-deployment auditing that monitors performance drift across sub-populations. Crucially, fairness must be framed in *clinical* rather than purely statistical terms—for instance, ensuring that false-negative rates do not delay treatment for already underserved groups.

#### 4.5 Privacy, consent and governance

Audio diaries, selfie videos and wearable biosignals are deeply personal. Even anonymised representations can be re-identified when modalities are combined. Patients may consent to one use (e.g. symptom tracking) but not anticipate secondary uses such as insurance underwriting. Robust governance therefore extends beyond encryption to include purpose-limiting licences, transparency dashboards that show how data are processed and revocable consent mechanisms. Differential privacy and homomorphic encryption offer technical safeguards, but their adoption hinges on computational feasibility and regulators' willingness to accept probabilistic rather than absolute anonymity.

#### 4.6 Evaluation in context

Traditional metrics such as accuracy or  $F_1$  flatten the nuanced objectives of mental-health care, where false negatives may delay lifesaving intervention and false positives can stigmatise. Models should therefore be judged in the context of the clinical pathway: How early is an impending relapse detected? Does the system triage scarce therapist time more effectively than standard practice? Randomised controlled trials, simulation studies of clinical workflows and decision-curve analyses that weigh benefit against harm are needed before deployment claims can be taken seriously.

#### 4.7 Complexity and comorbidity

Depression rarely appears in isolation; it co-occurs with anxiety, substance use or chronic pain. Multimodal models trained on single-label corpora can confuse overlapping phenotypes, leading to brittle predictions when presentations shift. Capturing this multifactorial reality demands hierarchical or multi-task formulations that share representations across disorders and account for temporal dynamics such as episode recurrence and treatment effects. Achieving this goal will require larger, longitudinal datasets annotated for multiple conditions and closer dialogue between ML researchers and clinical scientists.

#### 4.8 Toward ethically grounded, socially responsible solutions

Overcoming these challenges will require cross-sector collaboration: technologists to refine algorithms, clinicians to define clinically meaningful targets, ethicists to safeguard autonomy and justice, and patient advocates to ensure that lived experience shapes priorities. Only through such a holistic effort can multimodal ML mature from an intriguing research frontier into a dependable pillar of mental-health care.

### 5 Conclusion

Multimodal machine learning has shifted the conversation in digital psychiatry from *whether* machine intelligence can help clinicians to *how* we can deploy it safely and equitably. By fusing language, prosody, facial behaviour

and physiology, contemporary models already outperform unimodal baselines on screening, severity estimation and relapse prediction across depression, stress and related conditions.

For multimodal mental-health research to fulfil its promise, the field must address three intertwined imperatives at once. First, data ecosystems need to diversify and scale: longitudinal, culturally inclusive cohorts gathered under privacy-preserving and federated protocols are indispensable for learning representations that generalise across settings. Second, evaluation practice must evolve beyond leaderboards dominated by accuracy or F1; instead, benchmarks should probe cross-site transportability, calibration in the face of low prevalence, and—crucially—downstream effects on clinical decision-making. Third, systems must be designed for trustworthiness from the ground up, embedding clinically informed priors, offering faithful explanations, and undergoing rigorous fairness audits so they can satisfy regulators and earn patient acceptance. If these challenges are met, multimodal ML can evolve into a dependable pillar of mental-health care—supporting earlier detection, personalised interventions and continuous monitoring at a scale traditional services cannot match. Realising that future will require sustained collaboration among computer scientists, clinicians, ethicists and, critically, people with lived experience of mental illness. With such a coalition the field can move beyond proof-of-concept studies toward systems that measurably reduce suffering and broaden access to high-quality mental-health support worldwide.

## References

- [1] World Health Organization. Mental health, n.d. Accessed 28-Apr-2025.
- [2] World Health Organization. Mental health atlas 2020, 2020. Accessed 28-Apr-2025.
- [3] Mental health apps market size & share report (2030). <https://www.grandviewresearch.com/industry-analysis/mental-health-apps-market-report>, n.d. Grand View Research, accessed 28-Apr-2025.
- [4] A. Ahmad, V. Singh, and K. Upreti. A systematic study on unimodal and multimodal human-computer interface for emotion recognition. In *Computing, Internet of Things and Data Analytics (ICCI DA 2023)*, volume 1145 of *Studies in Computational Intelligence*. Springer Cham, 2024.
- [5] J. Gratch, R. Artstein, G. Lucas, G. Stratou, S. Scherer, A. Nazarian, F. Ringeval, B. Schuller, M. Valstar, N. Cummins, R. Cowie, L. Tavabi, and M. Schmitt. Avec 2019 workshop and challenge: State-of-mind, detecting depression with ai, and cross-cultural affect recognition. In *Proceedings of the 9th International Audio/Visual Emotion Challenge (AVEC '19)*, pages 3–12. ACM, 2019.
- [6] M. Valstar, B. Schuller, K. Smith, F. Eyben, B. Jiang, S. Bilakhia, S. Schnieder, R. Cowie, and M. Pantic. Avec 2013: The continuous audio/visual emotion & depression recognition challenge. In *Proceedings of the 3rd ACM International Workshop on Audio/Visual Emotion Challenge*, pages 3–10. ACM, 2013.
- [7] H. Cai, Z. Yuan, Y. Gao, et al. A multi-modal open dataset for mental-disorder analysis. *Scientific Data*, 9:178, 2022.
- [8] S. Koldijk, M. Sappelli, S. Verberne, M. Neerincx, and W. Kraaij. The swell knowledge work dataset for stress and user modelling research. In *Proceedings of the 16th ACM International Conference on Multimodal Interaction (ICMI 2014)*, Istanbul, Turkey, 2014.
- [9] J. Yoon, C. Kang, S. Kim, and J. Han. D-vlog: Multimodal vlog dataset for depression detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 12226–12234, 2022.
- [10] B. Zou. Chinese multimodal depression corpus. IEEE Dataport, 2022.
- [11] F. Ringeval, B. Schuller, M. Valstar, R. Cowie, H. Kaya, M. Schmitt, S. Amiriparian, N. Cummins, et al. Avec 2018 workshop and challenge: Bipolar disorder and cross-cultural affect recognition. In *Proceedings of the Audio/Visual Emotion Challenge and Workshop (AVEC '18)*, 2018.
- [12] Elvan Çiftçi, Heysem Kaya, Hüseyin Güleç, and Albert Ali Salah. The turkish audio-visual bipolar disorder corpus. In *2018 First Asian Conference on Affective Computing and Intelligent Interaction (ACII Asia)*, pages 1–6, 2018.
- [13] M. Sawadogo, F. Pala, G. Singh, I. Selmi, P. Puteaux, and A. Othmani. Ptsd in the wild: A video database for studying post-traumatic stress disorder recognition in unconstrained environments. *arXiv preprint*, 2022.
- [14] K. Kutt, D. Dražyk, L. Żuchowska, et al. Biraffe2: A multimodal dataset for emotion-based personalisation in rich affective game environments. *Scientific Data*, 9:274, 2022.
- [15] S. Koelstra, C. Mühl, M. Soleymani, S. Lee, J. et al. Deap: A database for emotion analysis using physiological signals. *IEEE Transactions on Affective Computing*, 3(1):18–31, 2012.
- [16] Mimansa Jaiswal, Cristian-Paul Bara, Yuanhang Luo, Mihai Burzo, Rada Mihalcea, and Emily Mower Provost. MuSE: a multimodal dataset of stressed emotion. In Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Marian, Hélène Mazo, Asuncion Moreno, Jan Odijk, and Stelios Piperidis, editors,

- Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 1499–1510, Marseille, France, May 2020. European Language Resources Association.
- [17] Renato Panda, Ricardo Malheiro, Bruno Miguel Machado Rocha, A. Oliveira, and Rui Pedro Paiva. Multi-modal music emotion recognition: A new dataset, methodology and comparative analysis. 2013.
  - [18] Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. MELD: A multimodal multi-party dataset for emotion recognition in conversations. In Anna Korhonen, David Traum, and Lluís Márquez, editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 527–536, Florence, Italy, July 2019. Association for Computational Linguistics.
  - [19] Jose A Miranda Calero, Laura Gutiérrez-Martín, Esther Rituerto-González, Elena Romero-Perales, Jose M Lanza-Gutiérrez, Carmen Peláez-Moreno, and Celia López-Ongil. Wemac: Women and emotion multi-modal affective computing dataset. *Scientific data*, 11(1):1182, 2024.
  - [20] Behnaz Nojavanasghari, Tadas Baltrušaitis, Charles E. Hughes, and Louis-Philippe Morency. Emoreact: a multimodal approach and dataset for recognizing emotional responses in children. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, ICMI ’16, page 137–144, New York, NY, USA, 2016. Association for Computing Machinery.
  - [21] P. Schmidt, A. Reiss, C. Duerichen, C. Marberger, and K. Van Laerhoven. Introducing wesad, a multimodal dataset for wearable stress and affect detection. In *Proceedings of ICMI 2018*, Boulder, USA, 2018.
  - [22] Seyedmajid Hosseini, Raju Gottumukkala, Satya Katragadda, Ravi Teja Bhupatiraju, Ziad Ashkar, Christoph W. Borst, and Kenneth Cochran. A multimodal sensor dataset for continuous stress detection of nurses in a hospital. *Scientific Data*, 9(1):255, June 2022.
  - [23] L. Stappen, A. Baird, L. Christ, M. Meßner, E. and B. Schuller. Muse-stress: Multimodal emotional stress (muse 2021) (version 1), 2021.
  - [24] K. Pisanski et al. Dataset for “multimodal stress detection: Testing for covariation in vocal, hormonal and physiological responses to trier social stress test”, 2018.
  - [25] M. Hosseini, F. Sohrab, R. Gottumukkala, et al. Empathieschool: A multimodal dataset for real-time facial expressions and physiological data analysis under different stress conditions. *arXiv preprint*, 2022.
  - [26] W. Lin, I. Orton, M. Liu, and M. Mahmoud. Automatic detection of self-adaptors for psychological distress. In *Proceedings of the 15th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2020)*, pages 371–378, 2020.
  - [27] V. Markova. Clas: A database for cognitive load, affect and stress recognition, 2020.
  - [28] L. Christ, S. Amiriparian, A. Baird, et al. The muse 2022 multimodal sentiment analysis challenge: Humour, emotional reactions and stress. In *Proceedings of MuSe 2022*, 2022.
  - [29] R. Meziati, Y. Beneszeth, P. De Oliveira, et al. Ubfcc–phys. IEEE Dataport, 2021.
  - [30] P. Mann, A. Paes, and H. Matsushima, E. See and read: Detecting depression symptoms in higher education students using multimodal social media data. *arXiv preprint*, 2019.
  - [31] T. Gui, L. Zhu, Q. Zhang, M. Peng, X. Zhou, K. Ding, and Z. Chen. Cooperative multimodal approach to depression detection in twitter. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 110–117, 2019.
  - [32] X. Zhang, B. Li, G. Qi, et al. A novel multimodal depression diagnosis approach utilising a new hybrid fusion method. *Biomedical Signal Processing and Control*, 96:106552, 2024.
  - [33] M. Rohanian, J. Hough, and M. Purver. Detecting depression with word-level multimodal fusion. In *Proceedings of Interspeech 2019*, pages 1443–1447, 2019.
  - [34] S. Liu, S. Wang, C. Sun, B. Li, S. Wang, and F. Li. Deepgcn based on variable multi-graph and multimodal data for asd diagnosis. *CAAI Transactions on Intelligence Technology*, 2024.
  - [35] Arnab Kumar Das and Ruchira Naskar. A deep learning model for depression detection based on mfcc and cnn generated spectrogram features. *Biomedical Signal Processing and Control*, 90:105898, 2024.
  - [36] A. Pradhan and S. Srivastava. Hybrid densenet with long short-term memory model for multi-modal emotion recognition from physiological signals. *Multimedia Tools and Applications*, 83:35221–35251, 2024.
  - [37] Z. Zhang, S. Zhang, D. Ni, Z. Wei, K. Yang, S. Jin, G. Huang, Z. Liang, L. Zhang, and L. Li. Multimodal sensing for depression risk detection: Integrating audio, video and text data. *Sensors*, 24(12):3714, 2024.
  - [38] K. Lilhore, U. S. Dalal, N. Varshney, K. Sharma, Y. B. Rao, K. V. Rao, V. R. Alroobaee, S. Simaiya, M. Margala, and P. Chakrabarti. Prevalence and risk-factors analysis of postpartum depression at early stage using a hybrid deep-learning model. *Scientific Reports*, 14:1–24, 2024.
  - [39] R. Beniwal and P. Saraswat. A hybrid bert–cnn approach for depression detection on social media using multimodal data. *The Computer Journal*, 67(7):2453–2472, 2024.
  - [40] S. Lee, Y. Cho, Y. Ji, M. Jeon, A. Kim, B. Ham, and Y. Joo, Y. Multimodal integration of neuroimaging and genetic data for diagnosis of mood disorders based on computer vision models. *Journal of Psychiatric Research*, 172:144–155, 2024.
  - [41] Y. Yao, M. Papakostas, M. Burzo, M. Abouelenien, and R. Mihalcea. Muser: Multimodal stress detection using emotion recognition as an auxiliary task. *arXiv preprint*, 2021.

- [42] Joseph Aina, Oluwatunmise Akinniyi, Md. Mahmudur Rahman, Valerie Odero-Marah, and Fahmi Khalifa. A hybrid learning-architecture for mental disorder detection using emotion recognition. *IEEE Access*, 12:91410–91425, 2024.
- [43] Ana-Maria Bucur, Adrian Cosma, Paolo Rosso, and Liviu P Dinu. It’s just a matter of time: Detecting depression with time-enriched multimodal transformers. In *European conference on information retrieval*, pages 200–215. Springer, 2023.
- [44] B. Yang et al. Mmpf: Multimodal purification fusion for automatic depression detection. *IEEE Transactions on Computational Social Systems*, 2024.
- [45] Y. Guo, C. Zhu, S. Hao, and R. Hong. A topic–attentive transformer–based model for multimodal depression detection. *arXiv preprint*, 2022.
- [46] Y. Tao, M. Yang, H. Li, Y. Wu, and B. Hu. Depmstat: Multimodal spatio–temporal attentional transformer for depression detection. *IEEE Transactions on Knowledge and Data Engineering*, 36(7):2956–2966, 2024.
- [47] A. Malhotra and R. Jindal. Xai transformer–based approach for interpreting depressed and suicidal user behaviour on online social networks. *Cognitive Systems Research*, 84:101186, 2024.
- [48] P. Moon and P. Bhattacharyya. We care: Multimodal depression detection and knowledge–infused mental–health therapeutic response generation. *arXiv preprint*, 2024.
- [49] Y. Xia, L. Liu, T. Dong, et al. A depression detection model based on multimodal graph neural network. *Multimedia Tools and Applications*, 83:63379–63395, 2024.
- [50] W. Zheng, L. Yan, C. Gou, and Y. Wang. Graph attention model embedded with multi–modal knowledge for depression detection. In *Proceedings of the IEEE International Conference on Multimedia and Expo (ICME 2020)*, pages 1–6, London, UK, 2020.
- [51] J. Cha, S. Kim, D. Kim, and E. Park. Mogam: A multimodal object–oriented graph attention model for depression detection. *arXiv preprint*, 2024.
- [52] H. Shen, S. Song, and H. Gunes. Multi–modal human behaviour graph representation learning for automatic depression assessment. *Apollo – University of Cambridge Repository*, 2024.
- [53] Abhinav Joshi, Ashwani Bhat, Ayush Jain, Atin Singh, and Ashutosh Modi. COGMEN: COntextualized GNN based multimodal emotion recognitioN. In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz, editors, *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4148–4164, Seattle, United States, July 2022. Association for Computational Linguistics.
- [54] Rina Carines Cabral, Soyeon Caren Han, Josiah Poon, and Goran Nenadic. Mm-emog: Multi-label emotion graph representation for mental health classification on social media. *Robotics*, 13(3), 2024.
- [55] F. Noman, M. Ting, C. H. Kang, R. Phan, and H. Ombao. Graph autoencoders for embedding learning in brain networks and major depressive disorder identification. *IEEE Journal of Biomedical and Health Informatics*, 28(3):1644–1655, 2024.
- [56] J. Pan, H. Lin, Y. Dong, Y. Wang, and Y. Ji. Mamf–gen: Multi–scale adaptive multi–channel fusion deep graph convolutional network for predicting mental disorder. *Computer Biology and Medicine*, 148:105823, 2022.
- [57] J. Hu, Y. Liu, J. Zhao, and Q. Jin. Mmgcn: Multimodal fusion via deep graph convolution network for emotion recognition in conversation. *arXiv preprint*, 2021.