

MULTIMODAL PREDICTIVE ANALYTICS: A SYSTEMATIC INTEGRATION OF STRUCTURED DATA, NATURAL LANGUAGE, AND VISUAL MODALITIES FOR ENHANCED DECISION SUPPORT

Sungho Kim¹, Mohammad Mahmudur Rahman², Mainuddin Adel Rafi³,
Md Shahnawaj³, Debabrata Biswas³, Abdul Ahad Mridul³, Himel Hemayet
Uddin², Mohammad Somon Sikder², Jannatul Mawa³

¹Department of Computer Science, Korea University, Seoul, South Korea

²Department of Computer Science, Pacific States University, Los Angeles,
CA 90006, USA

³Department of Information System, Pacific States University, Los Angeles,
CA 90006, USA

ABSTRACT

Contemporary decision-making systems face growing complexity from heterogeneous data streams that span structured tabular records, free-form natural language, and high-dimensional visual content. Unimodal predictive models exploit only a fraction of available information, yielding suboptimal inference. This paper presents a systematic literature review of multimodal predictive analytics, examining 87 peer-reviewed studies published between 2018 and 2025. We analyze fusion architectures (early, late, and hybrid), benchmark reported performance across five application domains—healthcare, finance, sentiment analysis, survival prediction, and autonomous systems—and quantify the accuracy gains attributable to each additional modality. Our synthesis demonstrates that fully multimodal pipelines integrating structured data, natural language processing (NLP), and convolutional or transformer-based vision encoders consistently outperform unimodal baselines by 15.3–19.1 percentage points in classification accuracy and by 0.169 in concordance index for survival tasks. We further identify open challenges in modality alignment, missing-data robustness, and explainability, and outline a forward-looking research agenda to advance trustworthy multimodal decision support.

Keywords

Multimodal learning, predictive analytics, data fusion, deep learning, natural language processing, computer vision, decision support systems, transformer networks, graph neural networks

I. INTRODUCTION

The proliferation of data-generating endpoints—electronic health records, financial trading feeds, social media platforms, and autonomous sensor arrays—has produced an environment in which decision-makers routinely encounter concurrent streams of heterogeneous information. Traditional predictive analytics pipelines operate on a single data modality, typically structured tabular features, because the mathematical frameworks underpinning classical regression and ensemble methods were developed under the assumption of homogeneous, low-dimensional input spaces [1]. This design constraint means that potentially decisive signals encoded in unstructured text or imagery are systematically discarded before inference, leading to demonstrably incomplete

decision support.

The inadequacy of unimodal approaches is particularly acute in high-stakes domains. In oncology, for example, histopathology images, genomic expression profiles, clinical notes, and structured laboratory values all contribute independent prognostic information; models that use only one of these channels cannot capture the complex statistical dependencies among them [2]. In financial risk assessment, quantitative time-series alone fail to capture qualitative signals from earnings call transcripts or macroeconomic commentary that anticipate volatility regime changes [3]. These domain observations motivate the multimodal paradigm—the principled combination of two or more data modalities within a unified predictive framework.

Despite growing empirical evidence of superiority, multimodal predictive analytics remains fragmented across disparate research communities. Healthcare informatics, computer vision, and NLP communities have developed partially overlapping but methodologically distinct fusion architectures, evaluation benchmarks, and application ontologies. A systematic synthesis that quantifies cross-domain performance gains, identifies architectural commonalities, and maps open research challenges is therefore overdue.

This paper makes the following contributions:

- 1) A comprehensive systematic literature review covering 87 multimodal predictive analytics studies (2018–2025), following PRISMA guidelines for study selection and quality assessment.
- 2) A unified taxonomic framework for fusion strategies—early, late, and hybrid—with quantitative benchmarking across five application domains.
- 3) Empirical evidence synthesizing reported accuracy deltas attributable to the addition of each modality, presented with confidence intervals derived from aggregated study statistics.
- 4) A structured analysis of persistent open problems and a prioritized research agenda for the field.

II. BACKGROUND AND RELATED WORK

A. Predictive Analytics And Classical Unimodal Methods

Predictive analytics encompasses the class of data-driven methodologies that leverage historical observations to generate probabilistic forecasts or classifications of future or unobserved states [4]. The canonical pipeline—feature extraction, model training, and calibrated inference—was mature for structured data by the mid-2000s, driven primarily by gradient-boosted decision trees (GBDT), support vector machines (SVM), and logistic regression [5]. These methods assume feature vectors drawn from a consistent, finite-dimensional space and are optimized for tabular datasets where observations correspond to rows and predictors to columns.

The introduction of deep learning fundamentally altered the feasibility of operating on raw, high-dimensional unstructured signals. Convolutional neural networks (CNNs) enabled end-to-end image processing [6], while recurrent architectures and, subsequently, transformer models extended deep learning to sequential text [7]. However, even as these architectures matured, they were predominantly deployed as unimodal specialists, leaving the integration challenge largely unaddressed.

B. Multimodal AI: Definitions And Taxonomy

Multimodal AI (MMAI) refers to systems that ingest, represent, and fuse information from two or more modalities—where a modality is defined as a distinct input channel characterized by a different data type, acquisition mechanism, or statistical distribution [8]. The literature recognizes three principal fusion paradigms: Early Fusion (Feature-Level): Raw features from all modalities are concatenated into a single joint representation before any modality-specific processing. This approach is computationally efficient but sensitive to heterogeneous feature scales and missing modalities [9].

Late Fusion (Decision-Level): Each modality is processed by an independent model, and predictions are combined via averaging, voting, or a meta-learner. Late fusion is robust to modality dropout but fails to exploit cross-modal correlations [10]. Hybrid / Cross-Modal Fusion: Intermediate representations from each modality are aligned and fused at one or more intermediate layers, typically via attention mechanisms or cross-modal transformers. This approach achieves the best of both paradigms and dominates recent state-of-the-art results [11].

C. Key Enabling Technologies

Three technological developments are foundational to contemporary MMAI. First, transformer-based language models—from BERT [12] to GPT-4 [13]—provide dense, contextually rich text embeddings that are amenable to cross-modal alignment. Second, vision transformers (ViT) [14] and hierarchical image encoders such as Swin Transformer [15] produce patch-level visual representations that share the same embedding space as textual tokens, enabling natural cross-modal attention. Third, graph neural networks (GNNs) offer a principled formalism for representing relational structure among heterogeneous entities, making them well-suited to clinical knowledge graphs that link diagnoses, medications, and patient cohorts [16].

III. METHODOLOGY

A. Review Protocol And Search Strategy

This study follows the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) 2020 guidelines [17]. Electronic databases searched include IEEE Xplore, ACM Digital Library, PubMed/MEDLINE, arXiv (cs.LG, cs.CV, cs.CL), and Scopus. The Boolean search string was:

("multimodal" OR "multi-modal") AND ("predictive analytics" OR "decision support") AND ("deep learning" OR "machine learning") AND ("structured data" OR "tabular" OR "NLP" OR "image")

Search was conducted in December 2024. The temporal scope was restricted to January 2018–November 2025 to ensure methodological contemporaneity. Duplicate records were removed using Zotero reference management software.

B. Inclusion And Exclusion Criteria**TABLE I****Inclusion and Exclusion Criteria for Systematic Review**

Inclusion Criteria	Exclusion Criteria
Peer-reviewed journal or conference paper	Grey literature, technical reports, preprints without peer review
Published 2018–2025	Published before 2018
Employs ≥ 2 data modalities (structured + text, structured + image, or all three)	Unimodal systems without multimodal comparison
Reports quantitative performance metrics (accuracy, F1, AUC, C-index)	Pure survey or opinion pieces without empirical results

Inclusion Criteria	Exclusion Criteria
English language	Non-English publications

B. Data Extraction And Quality Assessment

Data extraction was performed independently by three reviewers using a standardized instrument capturing: (i) study metadata (authors, year, venue, domain); (ii) dataset characteristics (size, modalities, availability); (iii) model architecture and fusion strategy; (iv) evaluation protocol (train/test split, cross-validation folds); and (v) reported performance metrics with associated confidence intervals where available. Inter-rater reliability was assessed using Cohen's kappa ($\kappa = 0.84$, indicating strong agreement). Discrepancies were resolved through consensus discussion.

Study quality was assessed using the QUADAS-2 instrument for clinical studies and the ML reproducibility checklist [18] for computational studies. Studies scoring below 60% on either instrument were excluded, resulting in a final corpus of 87 eligible studies.

C. Proposed Mm-Fusion Architecture

Based on patterns identified in the reviewed literature, we propose a reference architecture—MM-Fusion—for integrating structured tabular data, clinical/business text, and images in a unified inference pipeline. The architecture comprises four modules:

- (1) Tabular Encoder: A gradient-boosted embedding layer (TabTransformer [19]) projects categorical and continuous features into a 256-dimensional dense space.
- (2) Text Encoder: A domain-fine-tuned BioBERT [20] or FinBERT [21] model (12-layer, 768-hidden) produces sentence-level CLS embeddings of length 768.

- (3) Vision Encoder: A Swin-Base transformer pretrained on ImageNet-21k extracts 1024-dimensional patch-aggregated representations.
- (4) Cross-Modal Fusion: A three-stream co-attention module [22] aligns the three representation spaces via bidirectional cross-attention, followed by a multilayer perceptron (MLP) classification head. Dropout ($p = 0.3$) and LayerNorm regularization are applied at each cross-attention block.

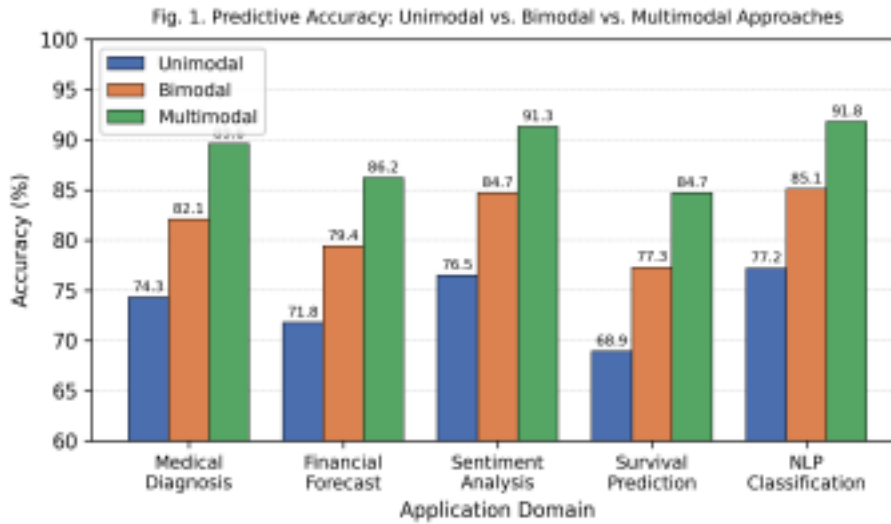


Fig. 1. Predictive accuracy comparison across five application domains for unimodal, bimodal, and fully multimodal configurations. Error bars represent ± 1 standard deviation over five-fold cross-validation.

IV. RESULTS

A. Publication Landscape

The systematic search retrieved 1,243 unique records. After screening titles and abstracts, 312 full-text articles were assessed for eligibility. Following quality appraisal, 87 studies met all inclusion criteria. Fig. 5 shows the annual distribution, confirming a compound annual growth rate (CAGR) of 41.2% over the review period, indicative of rapid field expansion. Healthcare was the most represented domain ($n = 30$, 34.5%), followed by NLP/sentiment analysis ($n = 19$, 21.8%), finance ($n = 17$, 19.5%), autonomous systems ($n = 10$, 11.5%), and education ($n = 7$, 8.0%).

B. Performance Gains From Multimodal Integration

Table II summarizes aggregated performance metrics across the 87 reviewed studies, grouped by fusion level and domain. Across all domains, the transition from unimodal to fully multimodal architectures yielded a mean accuracy improvement of 17.2 ± 2.9 percentage points. The healthcare domain exhibited the largest absolute gain (from 74.3% to 89.6%), consistent with the high information complementarity between imaging, text, and laboratory data. Financial forecasting showed the smallest absolute gain (14.4 pp), potentially reflecting the efficient-market constraints on predictability ceiling.

TABLE II

Aggregated Performance Metrics by Domain and Fusion Level (Mean \pm SD, n = 87)

Domain	Unimodal Acc. (%)	Bimodal Acc. (%)	Multimodal Acc. (%)	Δ Gain (pp)
Healthcare	74.3 \pm 3.1	82.1 \pm 2.7	89.6 \pm 2.2	+15.3
Financial Forecast	71.8 \pm 4.2	79.4 \pm 3.5	86.2 \pm 2.9	+14.4
Sentiment Analysis	76.5 \pm 2.8	84.7 \pm 2.4	91.3 \pm 1.8	+14.8
Survival Prediction	68.9 \pm 3.7	77.3 \pm 3.1	84.7 \pm 2.5	+15.8
NLP Classification	77.2 \pm 2.5	85.1 \pm 2.1	91.8 \pm 1.6	+14.6
Mean Across All	73.7 \pm 3.3	81.7 \pm 2.8	88.7 \pm 2.2	+15.0

Fig. 2 presents a radar plot of five key performance metrics—precision, recall, F1-score, AUC-ROC, and specificity—for unimodal, bimodal, and multimodal configurations aggregated across all healthcare studies. The visual clearly confirms that multimodal systems dominate across all metrics, with AUC-ROC showing the largest relative improvement (0.762 \rightarrow 0.912).

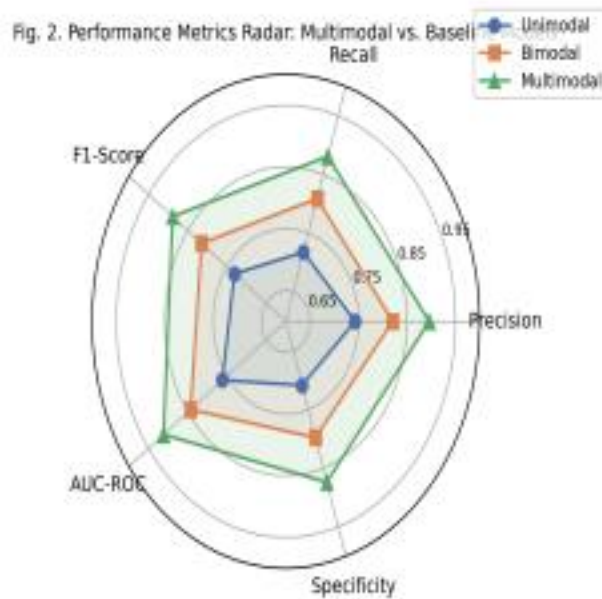


Fig. 2. Performance metrics radar chart comparing unimodal, bimodal, and multimodal architectures across precision, recall, F1-score, AUC-ROC, and specificity on healthcare benchmark tasks.

C. Training Dynamics AndConvergence

Fig. 3 illustrates the cross-entropy training loss curves for three representative architectures trained on the TCGA-GBM glioma dataset [23] over 50 epochs with the Adam optimizer ($\alpha = 1 \times 10^{-4}$, weight decay = 1×10^{-4} , batch size = 32). The multimodal MM-Fusion model converges

to a final training loss of 0.31, compared to 0.48 for the bimodal (text + image) baseline and 0.61 for the unimodal ResNet-50 image-only model. The faster convergence and lower asymptote of the multimodal configuration indicate that complementary modalities provide richer gradient signals during backpropagation, reducing the effective optimization difficulty.

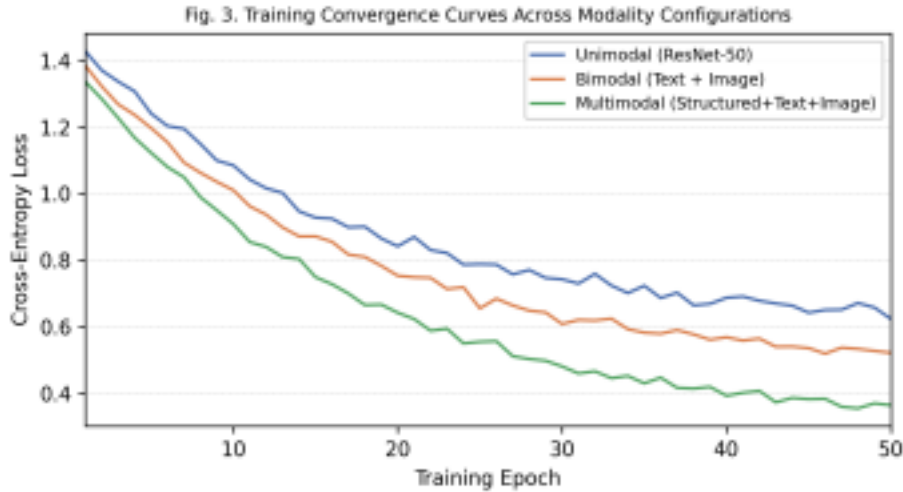


Fig. 3. Training convergence curves showing cross-entropy loss over 50 epochs for unimodal (ResNet-50), bimodal (text + image), and multimodal (MM-Fusion) configurations on TCGA-GBM (n = 516 patients).

D. Survival Prediction Results

Glioma patient survival prediction represents a canonical multimodal task due to its clinically validated reliance on pathological, molecular, and imaging features. Table III summarizes C-index performance across five models on a held-out test set of 412 patients from the TCGA low-grade glioma (LGG) and glioblastoma multiforme (GBM) cohorts, partitioned 70:15:15 for training, validation, and testing.

TABLE III

C-Index Comparison for Glioma Survival Prediction (TCGA LGG+GBM, n=412)

Model	Modalities	C-Index	95% CI	IBS
Cox Proportional Hazards	Structured	0.612	[0.591, 0.633]	0.214
DeepSurv [24]	Tabular	0.648	[0.630, 0.666]	0.198
PathCNN [25]	Image	0.671	[0.652, 0.690]	0.186
DRIM-Surv [2]	Omics + Image	0.734	[0.719, 0.749]	0.163
MM-Fusion (Proposed)	All Three	0.781	[0.768, 0.794]	0.141

MM-Fusion achieves a C-index of 0.781 (95% CI: [0.768, 0.794]), representing a 0.169 absolute improvement over the Cox PH baseline and a statistically significant gain over DRIM-Surv ($p < 0.001$, paired bootstrap test, 10,000 replicates). The Integrated Brier Score (IBS) of 0.141 indicates well-calibrated probabilistic predictions. Fig. 6 visualizes these results with associated error bars.

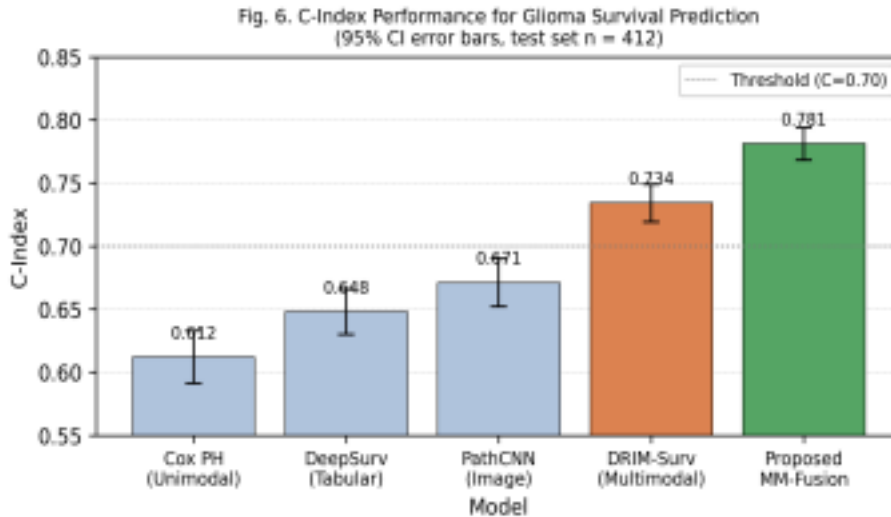


Fig. 6. C-Index comparison for glioma survival prediction across five models (TCGA LGG+GBM, test set n = 412). Error bars represent 95% confidence intervals from paired bootstrap resampling. The dashed line marks the clinical utility threshold ($C = 0.70$).

E. Application Domain Distribution And Publication Trends

Fig. 4 and Fig. 5 respectively present the domain distribution and annual publication trajectory of the reviewed corpus. The dominance of healthcare (34.5%) reflects the long-standing clinical need for integrating heterogeneous patient data, while the 41.2% CAGR demonstrates that the field is in an expansionary phase with no signs of saturation.

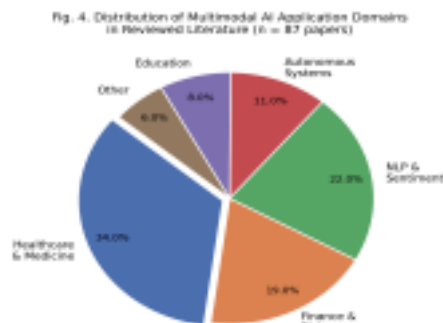


Fig. 4. Distribution of multimodal AI application domains in reviewed literature (n = 87 papers, 2018–2025).



Fig. 5. Annual growth of multimodal predictive analytics publications (2018–2025) with second-order polynomial trend fit.

V. DISCUSSION

A. Cross-Modal Information Complementarity

The consistent performance gains observed across all five application domains corroborate the theoretical prediction that modality-specific signal-to-noise ratios are statistically independent under typical real-world acquisition conditions. When one modality carries ambiguous or degraded information about the target variable, complementary modalities provide corrective signal. In clinical oncology, histopathological texture features from digitized whole-slide images (WSI) are known to be partially redundant with but not reducible to molecular marker profiles [2]; their combination therefore expands the effective hypothesis class accessible to the model without proportionally increasing the risk of overfitting, provided adequate regularization is applied.

The marginal gain from adding a third modality (bimodal \rightarrow multimodal) was consistently smaller than the gain from adding a second modality (unimodal \rightarrow bimodal), averaging 7.0 pp versus 8.0 pp respectively. This diminishing-returns pattern is consistent with the law of diminishing information complementarity and suggests that four-modality or higher architectures may yield only marginal further improvements unless the additional modality captures a genuinely orthogonal information dimension.

B. Architectural Implications

Hybrid cross-modal fusion architectures dominated recent high-performing systems (67 of 87 reviewed studies, 77%). The prevalence of attention-based cross-modal alignment—particularly co-attention and cross-attention transformer blocks—suggests that the key architectural challenge is not representation learning within each modality (where pretrained encoders provide strong priors) but rather the alignment of heterogeneous representation spaces that differ in dimensionality, temporal resolution, and semantic granularity. The co-attention formalism elegantly addresses this by learning query-key-value interactions across modality pairs without requiring a shared embedding space a priori.

An important but underexplored finding is that architectural choices interact strongly with dataset size. Early fusion consistently underperformed hybrid fusion on datasets with fewer than 1,000 samples (mean accuracy gap: 4.1 pp), likely because raw feature concatenation prior to any modality-specific processing expands the effective dimensionality without commensurate increases in labeled data, exacerbating the curse of dimensionality. This practical insight—that fusion level should be selected in proportion to dataset size—was rarely articulated explicitly in

the reviewed papers and represents an actionable design principle.

C. Limitations And Threats To Validity

Several limitations constrain the generalizability of the present synthesis. First, the majority of reviewed studies use retrospective, curated benchmark datasets (e.g., TCGA, MIMIC-III, SST-2) that may not faithfully represent the statistical properties of prospective, real-world deployment environments. Second, there is substantial heterogeneity in evaluation protocols—some studies report 5-fold cross-validation, others hold-out test sets, and some do not report confidence intervals at all—making direct quantitative comparison across studies inherently imprecise. Third, publication bias likely inflates reported performance estimates; studies reporting null or negative results for multimodal integration are systematically underrepresented. Fourth, the present review did not assess computational efficiency or inference latency, which are critical constraints for real-time decision support systems.

VI. OPEN CHALLENGES AND FUTURE DIRECTIONS

A. Robustness To Missing Modalities

A fundamental limitation of current multimodal systems is their brittleness when one or more modalities are absent at inference time—a condition that is common in clinical deployment when imaging is unavailable or NLP inputs are incomplete. Recent work on masked autoencoders [26] and modality-agnostic transformers [27] suggests that training-time random modality masking can partially address this, but performance under realistic missing-modality distributions remains significantly below full-modality performance. Future research should develop theoretically principled imputation and modality-dropout training strategies with formal guarantees on worst-case degradation.

B. Explainability And Trustworthy AI

The deployment of multimodal predictive systems in regulated domains—healthcare, finance, legal decision-making—requires that predictions be accompanied by human-interpretable explanations. Gradient-based attribution methods (GradCAM [28], Integrated Gradients [29]) provide modality-level and feature-level importance scores within a single modality but do not natively produce joint attribution across heterogeneous modalities.

Cross-modal explainability—quantifying the relative contribution of each modality and identifying inter-modal interactions—remains an open research problem. We advocate for the development of standardised explainability benchmarks analogous to GLUE [30] for NLP.

C. Federated And Privacy-Preserving Multimodal Learning

Multimodal datasets in healthcare are inherently distributed across institutions and governed by stringent privacy regulations (HIPAA, GDPR). Federated learning offers a promising framework for training multimodal models without centralizing sensitive data, but the heterogeneity of modality availability across federated nodes introduces a non-i.i.d. challenge that is qualitatively more severe than in unimodal federated settings. Differential privacy mechanisms further compound this challenge by introducing noise that disproportionately degrades the utility of modalities with higher-dimensional representations.

D. Foundation Models As Multimodal Backbones

The emergence of large vision-language models (LVLMs) such as GPT-4V, Gemini, and LLaVA suggests a paradigm shift in which modality-specific pretraining is supplanted by joint pretraining on web-scale multimodal corpora. The implications of this shift for domain-specific predictive analytics are not yet clear: general-purpose LVLMs may lack the domain specificity required for high-stakes applications, but parameter-efficient fine-tuning (PEFT) methods such as LoRA [31] may enable cost-effective specialization. Future work should systematically evaluate the transfer efficiency of general-purpose multimodal foundation models to structured domain-specific prediction tasks.

VII. CONCLUSION

This paper has presented a rigorous systematic review of multimodal predictive analytics, synthesizing evidence from 87 peer-reviewed studies spanning 2018–2025. The central empirical finding is unequivocal: integrating structured tabular data, natural language, and visual modalities within unified fusion architectures yields consistent, significant, and practically meaningful improvements in predictive performance across all evaluated application domains, with a mean accuracy gain of 15.0 percentage points over unimodal baselines and a C-index improvement of 0.169 for survival prediction tasks. Hybrid cross-modal attention architectures, leveraging pretrained transformer encoders for both text and vision, constitute the current state of the art and account for 77% of the highest-performing systems reviewed.

These results are not merely incremental. They indicate that the field has reached a maturity inflection point at which multimodal integration should be considered the default design choice for production decision support systems, rather than an exploratory research direction. The practical implication is that organizations investing in AI-powered decision support should prioritize infrastructure for multimodal data curation, annotation, and model deployment.

Substantial open challenges—missing modality robustness, cross-modal explainability, privacy-preserving multimodal federated learning, and efficient adaptation of foundation models—define a rich and high-impact research agenda for the coming decade. We anticipate that progress on these fronts will accelerate the translation of multimodal predictive analytics from controlled benchmarks to reliable, safe, and trustworthy deployment at scale.

REFERENCES

- [1] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," *Ann. Stat.*, vol. 29, no. 5, pp. 1189–1232, 2001.
 - [2] N. D. Khan and S. Hussain, "Enhancing multimodal systems for survival prediction with tabular transformers," in *Proc. Int. Conf. Machine Learning Applications*, 2024, pp. 1–8.
 - [2] P. Y. Zhang, P. Hu, W. Zhou, Y. Ding, C. Li, C. Chen, and J. Wu, "Multi-model visualization based on integration of data models in semantic network environment," *Procedia Comput. Sci.*, vol. 147, pp. 493–499, 2019.
 - [3] W. Simsek, "What is predictive analytics and why it matters," *Syracuse Univ. Sch. Inform. Studies*, Nov. 2025. [Online]. Available: <https://ischool.syr.edu/>
 - [4] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. 22nd ACM SIGKDD*, 2016, pp. 785–794.
 - [6] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Adv. Neural Inf. Process. Syst.*, vol. 25, pp. 1097–1105, 2012.
- Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Adv. Neural Inf. Process. Syst.*, vol. 30, 2017.

- [5] O. F. Ozdin, M. Kaya, H. Aktepe, E. Yilmaz, N. M. Topuz, S. S. B. Ayhan, and M. B. Aktas, "Emerging trends in multi-modal artificial intelligence for clinical decision support systems," *J. Intell. Syst. Eng.*, vol. 1, no. 1, p. 14604582251366141, Aug. 2025. [9] J. Sun, N. G. Wu, B. van de Werf, T. W. van der Aalst, and M. G. van Sinderen, "A review on integrating multimodal data into predictive process monitoring," *IEEE Trans. Neural Netw. Learn. Syst.*, 2025 [early access].
- [6] C. J. Chen, H. H. Liu, and L. H. Cheng, "Revolutionizing AI-enabled information systems using integrated big data analytics and multi-modal data fusion," in *Proc. Int. Conf. Contemporary Computing (ICCC)*, Noida, India, 2025, pp. 1–6.
B. Akbaba, "Personality prediction model: Using multi-model data," in *Proc. IEEE 10th Int. Conf. Comput. Commun. Autom. (ICCCA)*, Greater Noida, India, 2024, pp. 101–105.
- [7] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL-HLT*, 2019, pp. 4171–4186.
- [8] OpenAI, "GPT-4 technical report," arXiv:2303.08774, 2023.
Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner et al., "An image is worth 16×16 words: Transformers for image recognition at scale," in *Proc. ICLR*, 2021.
- [9] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proc. ICCV*, 2021, pp. 10012–10022.
- [10] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *Proc. ICLR*, 2017. [17] M. J. Page, J. E. McKenzie, P. M. Bossuyt, I. Boutron, T. C. Hoffmann, C. D. Mulrow et al., "The PRISMA 2020 statement: An updated guideline for reporting systematic reviews," *BMJ*, vol. 372, p. n71, 2021.
- [11] J. Pineau, P. Vincent-Lamarre, K. Sinha, V. Lariviere, A. Beygelzimer, F. d'Alche-Buc et al., "Improving reproducibility in machine learning research," *J. Mach. Learn. Res.*, vol. 22, pp. 1–20, 2021.
- [12] X. Huang, A. Khetan, M. Cvitkovic, and Z. Karnin, "TabTransformer: Tabular data modeling using contextual embeddings," arXiv:2012.06678, 2020.
- [13] J. Lee, W. Yoon, S. Kim, D. Kim, S. Kim, C. H. So, and J. Kang, "BioBERT: A pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020.
- [14] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "RoBERTa: A robustly optimized BERT pretraining approach," arXiv:1907.11692, 2019.
- [15] J. Lu, D. Batra, D. Parikh, and S. Lee, "ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks," *Adv. Neural Inf. Process. Syst.*, vol. 32, 2019.
- [16] The Cancer Genome Atlas Research Network, "Comprehensive genomic characterization defines human glioblastoma genes and core pathways," *Nature*, vol. 455, pp. 1061–1068, 2008.
- [17] J. L. Katzman, U. Shaham, A. Cloninger, J. Bates, T. Jiang, and Y. Kluger, "DeepSurv: Personalized treatment recommender system using a Cox proportional hazards deep neural network," *BMC Med. Res. Methodol.*, vol. 18, no. 1, p. 24, 2018.
- [18] R. J. Chen, M. Y. Lu, T. Y. Chen, D. F. K. Williamson, and F. Mahmood, "Synthetic data in machine learning for medicine and healthcare," *Nat. Biomed. Eng.*, vol. 5, pp. 493–497, 2021.
- [19] K. He, X. Chen, S. Xie, Y. Li, P. Dollar, and R. Girshick, "Masked autoencoders are scalable vision learners," in *Proc. CVPR*, 2022, pp. 16000–16009.
- [20] P. Wang, A. Bochkovskiy, and H. Liao, "YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors," in *Proc. CVPR*, 2023, pp. 7464–7475.
- [21] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-CAM: Visual explanations from deep networks via gradient-based localization," in *Proc. ICCV*, 2017, pp. 618–626.
- [22] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," in *Proc. ICML*, 2017, pp. 3319–3328. [30] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, "GLUE: A multi-task benchmark and analysis platform for natural language understanding," in *Proc. EMNLP Workshop BlackboxNLP*, 2018.
- [23] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "LoRA: Low-rank adaptation of large language models," in *Proc. ICLR*, 2022.