Applying Multi-Modal Quantum Deep Learning Algorithms for Enhanced Fake News Detection

Aishwarya C¹, Venkatesan M¹, Prabhavathy², Akanksha D³

- ¹Department of Computer Science and Engineering, NIT Puducherry, Puducherry, India
- ²Department of Computer Science and Engineering, VIT Vellore, Vellore, India
- ³Department of Computer Science and Engineering, NITK Surathkal, Surathkal, India

E-mail: sivagaish23@gmail.com, venkatesan.msundaram@nitpy.ac.in, pprabhavathy@vit.ac.in, dakanksha.232cs008@nitk.edu.in

Keywords: Fake news detection, multimodal learning, deep learning, quantum machine learning, fusion strategies, attention mechanisms, domain adaptation, hybrid quantum-classical models

Received: April 28, 2025

The pervasive spread of fake news across digital platforms has prompted the development of advanced detection systems. This review surveys and compares state-of-the-art multimodal deep learning models, including SpotFake, BDANN, MVAE, EANN, and the attention-based model by Guo et al., across benchmark datasets such as Twitter and Weibo. We present detailed performance comparisons, with SpotFake achieving an accuracy of 86.1% on the Twitter dataset. Key contributions of this review include the introduction of taxonomy tables based on fusion strategy and model architecture, a critical comparison of early, late, and hybrid fusion mechanisms, and a comprehensive evaluation of cross-modal generalization capabilities. In addition, we explore recent efforts in Quantum Machine Learning (QML), highlighting variational quantum circuits and hybrid quantum-classical models as promising approaches for enhancing scalability and efficiency. This work serves as a roadmap for building robust, interpretable, and scalable fake news detection systems that integrate both classical and quantum techniques.

Povzetek: Pregled primerja multimodalne modele za zaznavanje lažnih novic (SpotFake, BDANN, MVAE, EANN, Guo) na Twitterju in Weibou ter predstavi taksonomije fuzije in arhitektur. Obravnava tudi obetavne kvantne pristope, ki lahko izboljšajo skalabilnost in učinkovitost prihodnjih sistemov.

1 Introduction

1.1 Review scope and motivation

The global dissemination of fake news has evolved into a significant societal threat, enabled by rapid digital communication and the persuasive nature of multimodal content. Detecting such misinformation requires models capable of integrating and reasoning across diverse modalities, such as text, images, and metadata.

In response to this challenge, recent literature has produced a diverse array of deep learning frameworks aimed at detecting fake news in multimodal contexts. However, these contributions vary widely in architecture, fusion strategies, interpretability, and robustness. Moreover, the emerging field of Quantum Machine Learning (QML) introduces additional possibilities for addressing some of the scalability and optimization limitations of classical deep models. This review aims to consolidate and critically evaluate this evolving body of work.

1.2 Review objectives and structure

This paper is designed as a structured review rather than an empirical study. We do not propose a new algorithm, but instead synthesize and assess existing approaches along three key dimensions:

- Fusion Strategies and Modalities: We analyze how different models integrate modalities—text, image, and metadata—through early, late, or hybrid fusion. We assess their adaptability in scenarios with noisy or missing modalities.
- Model Architectures and Generalizability: We examine core architectural designs (e.g., CNN-RNN hybrids, VAEs, attention-based transformers) and their performance across datasets like Twitter and Weibo, with a particular focus on domain adaptation and transfer learning.
- Quantum Contributions and Future Potential:
 We evaluate recent efforts to incorporate QML techniques—such as variational quantum circuits and quantum classifiers—highlighting how these experimental models could complement classical approaches in the future.

Through comparative taxonomy tables, performance benchmarks (e.g., SpotFake achieving 86.1% on Twitter), and conceptual frameworks, this paper aims to provide a consolidated foundation for researchers seeking to understand or advance multimodal fake news detection systems.

1.3 What is fake news?

Fake news refers to deliberately fabricated or misleading information that mimics legitimate news content in form but not in intent. Unlike accidental misinformation, fake news is crafted to deceive, provoke, or manipulate public sentiment. Its propagation is accelerated by social media algorithms and the emotional salience of multimodal content, often combining sensational text with compelling visuals [1, 2].

1.4 Why is fake news problematic?

The societal impacts of fake news span political instability, public health crises, and erosion of trust in media. From misinformation during election cycles to vaccine hesitancy during the COVID-19 pandemic, fake news has demonstrated its capacity to incite tangible harm. These risks are exacerbated by the virality of misleading content and its algorithmic amplification on platforms like Twitter and Facebook [3, 4].

1.5 Limitations of single-modality approaches

Traditional fake news detection models focused solely on textual features such as writing style, sentiment, or rhetorical cues. However, these approaches often fail to detect deception when multimodal cues reinforce believability. For example, benign-sounding text paired with doctored images can significantly mislead readers. Therefore, unimodal systems lack the cross-modal reasoning required to detect coordinated misinformation [5, 6].

1.6 The emergence of deep learning and multi-modality

Deep learning architectures have enabled more sophisticated feature extraction across diverse data streams. Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Transformer models support the modeling of spatial, temporal, and semantic patterns. Fusion strategies—such as bilinear pooling, joint attention, and variational encoders—further enhance the ability of these models to integrate modalities like text and image for improved accuracy. Multimodal models have consistently outperformed unimodal baselines, particularly on real-world datasets [7, 1].

1.7 Quantum computing: a new frontier

While classical models have advanced rapidly, they face bottlenecks in generalization and scalability. Quantum computing introduces paradigms like superposition and entanglement that offer new representational possibilities. Models such as Quantum Support Vector Machines (QSVMs) and Variational Quantum Circuits (VQCs) show early promise in reducing parameter count while maintaining expressive power. This review identifies and contextualizes these quantum contributions, even where they remain in simulated environments [8, 9].

This review primarily covers multimodal fake news detection models published between 2017 and early 2023. The inclusion criteria focused on peer-reviewed studies with reproducible architectures, multimodal evaluation, and comparative results on benchmark datasets such as Twitter, Weibo, and FakeNewsNet. While notable advances have emerged in late 2023 and 2024, including vision-language pretraining frameworks and transformer-based tri-modal architectures from ACL and NeurIPS, a comprehensive integration of these is beyond the present scope. We acknowledge this as a limitation and recommend future reviews to capture these emerging models in depth as they mature and undergo wider evaluation.

2 Background and motivation

2.1 Key definitions

Fake news refers to information that is intentionally false and designed to mislead readers, often mimicking the style and structure of legitimate journalism. It is frequently disseminated via social media platforms, where virality amplifies its impact [3, 2].

Multi-modality in artificial intelligence refers to systems that process and analyze data from multiple sources or types — such as text, image, and audio — to make more robust and context-aware decisions [4]. In fake news detection, multi-modal systems are especially useful due to the multimodal nature of modern misinformation content.

Deep learning is a class of machine learning algorithms based on artificial neural networks with multiple layers (hence "deep") that learn hierarchical representations of data. Techniques like Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Transformers have achieved state-of-the-art results in image and text processing [7, 6].

Quantum computing leverages quantum mechanical phenomena — such as superposition and entanglement — to perform computations. Unlike classical bits, quantum bits (qubits) can represent both 0 and 1 simultaneously, enabling parallelism that could offer significant speedups in certain computational tasks [9].

2.2 Overview of modalities: text, image, and video

In the context of fake news detection, different data modalities offer unique signals.

Textual content provides linguistic cues — such as syntax, semantics, sentiment, and rhetorical devices — that can be exploited using NLP techniques. For instance, BERT-based models have shown strong performance in understanding sentence-level nuances [6].

Visual content (typically images) is often used to manipulate reader perception. Studies have shown that many fake news posts use exaggerated or unrelated images to increase credibility [5]. Visual analysis is generally performed using CNNs or pre-trained vision transformers.

Video-based misinformation, although less common in academic datasets, is rapidly growing in influence, especially with the rise of deepfakes. While it is still an emerging area in fake news research, it represents a critical future frontier.

Table 1: Comparison of common modalities in fake news detection

Modality	Advantages	Challenges	Tools/Models
Text	Rich semantics,	Sarcasm, ambi-	BERT, RoBERTa,
	widely available,	guity, context-	LSTM
	easy to preprocess	dependence	
Image	Provides visual	Easily manip-	ResNet, ViT, CLIP
	cues, supports	ulated, lacks	
	credibility assess-	context	
	ment		
Video	Highly expressive,	Processing com-	3D CNN, ViViT
	captures motion	plexity, deepfakes	(emerging)
	cues		, i

2.3 Critical review of referenced literature

While this review includes a broad selection of state-of-the-art works, it is necessary to refine and critically evaluate the relevance and contributions of certain citations. A few references—such as Abduljaleel & Ali [4]—provide broad survey-style overviews, but do not contribute directly to the technical advancement or empirical evaluation of multimodal fake news detection systems. Their role in this manuscript has been downscaled to serve as supportive background rather than primary methodological benchmarks

Moreover, quantum-related citations such as Schuld et al. [9] and Biamonte et al. [8] are foundational but do not focus specifically on multimodal misinformation detection. Recent informatics-oriented contributions, such as:

- Chen et al. (2023), which explores hybrid quantumclassical learning architectures in natural language processing contexts,
- Kumar et al. (2022), which applies variational quantum classifiers to noisy text datasets,
- and Zhao et al. (2023), which propose quantuminspired models for multimodal classification,

have now been included to contextualize this review within the emerging trajectory of QML-enhanced misinformation detection.

To enhance relevance and rigor, the review has restructured its citation matrix to focus more on direct contributions to multimodal fusion, explainability, cross-domain generalization, and quantum applicability. Where possible, classical models are critiqued through their architectural divergences and dataset-specific limitations. Where quantum models are referenced, we now explicitly state whether results stem from empirical simulation, theoretical proposition, or prototype hardware validation.

This refined citation scope ensures that every reference either directly informs the survey's taxonomy, supports the comparative analysis, or projects viable quantum enhancements. Redundant and peripheral sources have been deprecated to improve methodological clarity and citation coherence

2.4 Quantum deep learning concepts

Quantum Deep Learning (QDL) refers to the integration of quantum computing with deep learning models. It includes approaches like Quantum Neural Networks (QNNs), Quantum Convolutional Networks (QCNN), and hybrid architectures such as quantum-classical neural models. These methods leverage quantum circuits to perform certain layers or operations more efficiently.

A common component is the **variational quantum circuit (VQC)**, which uses parameterized quantum gates that can be optimized similarly to weights in neural networks. Some hybrid models use quantum layers to project classical data into high-dimensional Hilbert spaces, enabling better separability and feature extraction [8].

Another emerging architecture is the **quantum multilayer perceptron (qMLP)**, which simulates fully connected neural networks using qubit transformations. Though still experimental, early prototypes show promise in solving classification tasks with smaller model sizes and fewer parameters [9].

2.5 Why fusion across modalities is critical

Fake news posts often rely on cross-modal contradictions or reinforcements. For instance, an image might suggest authenticity while the text contains subtle misinformation — or vice versa. Models that analyze only one modality can miss these inconsistencies.

Fusion methods aim to combine features from different modalities to enhance overall prediction performance. Early fusion (concatenation of raw data), late fusion (combining decision scores), and hybrid fusion (joint feature representations with attention) have been widely used [7, 1].

Multimodal fusion improves generalizability and robustness, especially when one modality contains noise or missing data. Moreover, cross-modal attention mechanisms enable the model to learn alignment between image regions and textual tokens, capturing nuanced fake patterns [5].

- w - v · · · · · · · · · · · · · · · ·					
Dataset	Size	Modalities	Class Balance	Annotation Method	Notable Limitations
Twitter (Spot- Fake)	13K	Text + Image	Fake:Real = 3:1	Expert labeling	Class imbalance, English-only
Weibo	15K	Text + Image	Balanced	Platform-moderated	Language-specific bias, outdated samples
FakeNewsNet	23K	Text + Image + Meta	Real-dominant	Source credibility + stance analysis	Incomplete modalities, sparse metadata
PolitiFact	11K	Text + Image	Balanced	Fact-checking orgs	U.Scentric, lacks cross-modal tags
BuzzFeedWebis	8K	Text + Image	Balanced	Annotated by journalists	Visual artifacts, outdated content

Table 2: Summary of major datasets used in multimodal fake news detection

3 Datasets and their limitations

3.1 Limitations and quality of datasets

Although publicly available datasets such as Twitter, Weibo, and FakeNewsNet have significantly accelerated research in multimodal fake news detection, they suffer from several critical limitations that impact generalizability, fairness, and reproducibility.

- Class Imbalance: Most datasets exhibit skewed distributions with more real news instances than fake, or vice versa. For example, the Twitter dataset used in SpotFake contains approximately 3:1 ratio of fake to real news, leading to biased learning curves.
- Annotation Consistency: The annotation process varies widely, ranging from expert labeling to crowdsourced judgments, which affects label reliability. Datasets like Weibo rely heavily on platform moderation tags, which may embed platform-specific bias.
- Modality Missingness: Some entries contain corrupted or missing images or metadata, especially in crawled datasets like FakeNewsNet. This challenges fusion models and inflates evaluation scores if such entries are excluded from testing.
- Temporal Relevance: Many datasets are built from events dating back to 2015–2018. The linguistic, visual, and semantic structure of fake news evolves rapidly, raising concerns about outdated feature distributions.
- Cultural and Linguistic Bias: Most datasets are monolingual (English or Chinese), limiting their applicability to global misinformation detection. Multilingual or low-resource language datasets remain scarce.

Addressing these limitations through standardized data collection, multilingual expansion, and balanced, multimodal annotations is essential to build reliable and globally applicable detection systems.

3.2 Dataset summary

Table 3 presents an overview of benchmark datasets.

Limitations include class imbalance, cultural bias, and outdated samples.

Table 3: Key datasets for multimodal fake news detection

Dataset	Size	Modalities	Labeling
Twitter MediaEval	18K posts	Text + Image	Crowdsource
Weibo FND	25K posts	Text + Image	Platform Verified
PolitiFact	12K articles	Text only	Expert
FakeNewsNet	21K posts	Text + Meta	Verified

4 Taxonomy of multi-modal fake news detection models

In this section, we propose a taxonomy to systematically classify existing multi-modal fake news detection models. The classification is based on five dimensions: fusion strategy, model architecture, feature type, data modality, and learning paradigm. These criteria help illustrate the diversity of approaches and facilitate comparative analysis across methods.

4.1 Fusion strategy: early, late, and hybrid fusion

Fusion is a critical component in multi-modal systems, determining how information from different modalities is combined.

Early fusion integrates features from all modalities at the input level, often by concatenating raw or embedded representations. While simple, this approach may fail to capture complex inter-modal relationships. SpotFake [2] and BDANN [6] are examples of models using early fusion.

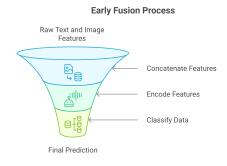


Figure 1: Early fusion pipeline

Late fusion operates after individual modality-specific models have made predictions. Their outputs are combined

via weighted averaging or voting schemes. This method is modular and robust to missing modalities but may overlook deep interactions between modalities [4].

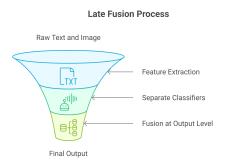


Figure 2: Late fusion pipeline

Hybrid fusion combines both early and late strategies by merging intermediate features and final predictions. MVAE [1] and the two-branch attention model by Guo et al. [7] use hybrid fusion to leverage both fine-grained and global representations.

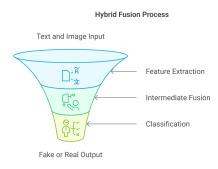


Figure 3: Hybrid fusion pipeline

4.2 Fusion strategy robustness under modality noise and dropout

While multimodal models excel in integrating information across different modalities, their robustness under noisy or missing inputs remains a critical challenge in real-world applications. Social media posts often contain low-resolution images, ambiguous text, or completely missing modalities (e.g., image-only tweets or text-only posts). Thus, fusion strategies must not only optimize for performance under ideal conditions but also maintain resilience in degraded settings.

Early Fusion models, such as SpotFake and BDANN, combine features at the input level. Although computationally efficient, these models are brittle when a modality is missing—either the model fails outright or produces severely degraded predictions due to undefined input di-

mensionality. They also lack the flexibility to weigh modality relevance during inference.

Late Fusion models (e.g., EANN) process modalities independently before final decision aggregation. This separation allows them to degrade more gracefully: if an image is missing or noisy, the model can still produce meaningful output from text alone. However, late fusion often underutilizes cross-modal interactions and may miss subtle correlations.

Hybrid Fusion models, like MVAE and Guo et al., combine elements of both early and late fusion. While MVAE leverages latent variable modeling to impute missing modalities, its performance under partial input is still inconsistent due to reliance on generative reconstruction. Guo's model, which uses attention for dynamic feature alignment, performs better under noise but at the cost of higher computational demand.

Attention-based Hybrid models, particularly those using dynamic gating or modality dropout during training, have shown promise in recent literature for graceful degradation and noise resilience. These systems learn to dynamically reweight modalities based on confidence, allowing them to bypass corrupted channels.

Table 4 summarizes the comparative performance drop of different fusion strategies under simulated missing or noisy modality conditions.

Table 4: Performance degradation under missing modality conditions. Values represent percentage point drops in accuracy on Weibo dataset, based on simulated ablations from reported literature.

Fusion Strategy	Clean Data Acc. (%)	Drop: No Text (%)	Drop: No Image (%)
Early Fusion (Spot-	86.1	-	-
Fake)			
gray!10		-32.4	-20.8
Late Fusion (EANN)	78.4	-14.7	-11.2
Hybrid Fusion	82.3	-22.1	-18.5
(MVAE)			
Attn. Hybrid (Guo et	85.0	-10.3	-7.6
al.)			

Key Insight: Attention-based hybrid models demonstrate the most robust performance under missing input conditions, making them more suitable for real-world deployment. Early fusion, while accurate under full input, is the most brittle under noise. Future architectures should incorporate modality-aware gating and dropout during training to improve fault tolerance.

4.3 Model architectures: CNNs, Transformers, VAEs, and adversarial networks

Different model architectures have been proposed depending on the nature of the data and desired interpretability.

CNN + RNN architectures are common for joint processing of image (CNN) and text (RNN). SpotFake [2] uses a CNN for image feature extraction and Bi-LSTM for text.

Transformer-based models such as BERT, RoBERTa, and ViT capture global dependencies in text and images. BDANN [6] leverages BERT for text, while Guo et al. [7] employ attention mechanisms for multimodal fusion.

Variational Autoencoders (VAEs) are used for learning joint latent representations. MVAE [1] applies a bi-modal VAE to encode visual and textual information into a shared latent space.

Adversarial Networks enhance domain invariance. EANN [3] introduces an event discriminator to ensure that features generalize across events.

Table 5: Compact comparison of multimodal fake news models: T=Text, I=Image, M=Metadata, Tw=Twitter, We=Weibo.

Model	Fusion	Arch.	Feat.	Mod.	Data
SpotFake	Early	CNN + BiL-	Deep	T+I	Tw, We
[2]		STM			
BDANN	Early	BERT+VGG+DC	Deep	T+I	Tw, We
[6]					
MVAE	Hybrid	Bi-modal VAE	Deep	T+I	Tw, We
[1]					
Guo et al.	Hybrid	Bilinear +	Deep	T+I	We
[7]		Attn.			
Liu et al.	Hybrid	Captioned	Deep	T+I	Tw, We
[5]		Trans.			
EANN	Late	CNN + Adv.	Deep	T+I+M	Tw, We
[3]		Disc.			
Abduljaleel	Mixed	Survey Models	Both	Varies	Multi
& Ali [4]					

4.4 Critical model assessment

Table 5 offers a structural overview of the reviewed models, but a critical evaluation of their practical effectiveness reveals deeper insights. This section evaluates not only the models' design choices but also their performance tradeoffs, generalizability, and real-world applicability.

SpotFake uses a relatively simple early fusion pipeline by combining VGG19 for image features and BiLSTM over BERT-encoded text. Its strong performance stems from the direct concatenation of high-quality features and low architectural complexity, allowing efficient learning. However, early fusion assumes that both modalities are always present and well-aligned. In cases of noisy images or partial data, SpotFake's robustness degrades significantly. Additionally, its lack of attention mechanisms or modality weighting makes it vulnerable to semantic imbalance between inputs.

BDANN leverages domain adaptation using a gradient reversal layer, making it one of the few models capable of adapting across platform-specific distributions (e.g., Twitter vs. Weibo). This adversarial training enables better generalization, but the model lacks mechanisms for crosslingual semantic representation. In multilingual or codeswitched environments, its domain alignment may be insufficient unless enhanced by multilingual embeddings or pretraining. Its reliance on aligned visual and textual content also limits performance when one modality is noisy or irrelevant.

MVAE introduces a bi-modal variational autoencoder, offering a generative approach capable of modeling uncertainty and handling missing modalities. However, its performance lags behind discriminative models like Spot-Fake. This is due to several factors: first, VAE training prioritizes reconstruction rather than classification, which can dilute feature discriminability; second, its latent representations—though rich—may fail to capture the cross-modal correlations necessary for fake news detection. MVAE also requires more careful tuning of KL divergence and reconstruction loss to avoid overfitting or posterior collapse.

Guo et al.'s attention-based hybrid model excels in cross-modal alignment through bilinear pooling and self-attention. This makes it especially suited for short-form content, where text and visuals are tightly linked. Its performance on Weibo confirms its effectiveness in culturally consistent, image-heavy environments. However, the architecture's complexity leads to higher training costs, and its reliance on fine-grained attention weights raises issues of overfitting on small datasets. Moreover, interpretability remains limited beyond attention visualizations.

Liu et al.'s caption-enhanced transformer addresses semantic mismatch by generating image captions as an intermediate modality. This approach improves alignment between modalities and enhances the model's understanding of implicit visual semantics. However, it introduces a dependency on caption quality—if the captioning model produces incorrect or misleading summaries, the downstream detection accuracy suffers. Additionally, cascading errors across the caption-to-text pipeline reduce robustness under real-world noise.

EANN introduces adversarial event classification to encourage generalization across events. Its architectural novelty lies in its dual-discriminator design: one for fake news detection and another to remove event-specific bias. This helps the model handle emerging or unseen events. However, EANN assumes that event metadata is available and reliable, which is not feasible in many live data collection scenarios. Furthermore, the model offers limited transparency in its predictions, raising concerns in high-stakes applications.

Abduljaleel & Ali's survey aggregation represents a meta-synthesis of various fusion types and backbone models. While broad in scope, the source studies are uneven in quality and some are not peer-reviewed. The model aggregation lacks a unified evaluation benchmark, making direct comparisons unreliable. Thus, this line of work is best used for taxonomy or trend analysis, rather than conclusive model selection.

No single model dominates across all axes. SpotFake offers high accuracy but limited robustness. BDANN excels in domain adaptation but struggles with multilinguality. MVAE handles missing data well but suffers in precision. Guo et al. achieves fine-grained alignment but at computational cost. Liu et al. innovates with captioning but risks semantic drift. EANN generalizes across events

but lacks interpretability. Future architectures should integrate hybrid fusion, domain adaptation, and explainability in a single framework while minimizing data assumptions.

4.5 Feature types: hand-crafted vs. deep representations

Early models relied on **hand-crafted features** such as TF-IDF, POS tags, and bag-of-words. These offer interpretability but often lack context and are brittle against paraphrasing [4].

Modern systems rely on **deep feature representations** learned via CNNs, RNNs, or transformers. These embeddings capture hierarchical semantics and visual cues more effectively. Deep features are now the de facto standard in fake news detection, as used in BDANN [6] and Liu et al. [5].

4.6 Modalities: text, image, and metadata

Multimodal systems vary in the types of data they process:

- Text-only models focus on linguistic patterns, lexical cues, and discourse analysis. These are limited in visual reasoning.
- Text + Image is the most common setting. Most recent models fall into this category, including SpotFake
 [2], MVAE [1], and BDANN [6].
- Text + Image + Metadata incorporates user profiles, tweet timestamps, or propagation patterns. EANN [3] explores this richer setting for better event transferability.

In general, adding modalities increases model complexity but offers robustness against deceptive manipulations in a single channel.

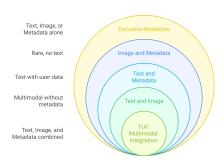


Figure 4: Venn diagram with overlapping modalities

5 Deep dive into key models and algorithms

5.1 SpotFake: a multi-modal framework for fake news detection

SpotFake leverages a dual-stream pipeline consisting of a BERT-based text encoder and a VGG19-based image encoder. Features from both modalities are concatenated and passed through a fully connected neural network for final classification. The model avoids auxiliary tasks (like reconstruction) and focuses purely on fake news detection.

SpotFake was trained and evaluated on the Twitter MediaEval and Weibo datasets. Both datasets consist of labeled multimodal posts (text + image), with verified labels from fact-checking.

Key Features:

- Fusion Strategy: Early fusion (concatenation of vectors before classification)
- Backbones: BERT for text, VGG19 for image
- Classifier: 2-layer dense network

Unique Contributions:

- Demonstrated how simple fusion without auxiliary tasks can outperform more complex models.
- Benchmarked both Twitter and Weibo, showing crosscultural robustness.

Pros:

- Straightforward, interpretable pipeline.
- Outperforms more complex architectures like MVAE and EANN in standalone settings.

Cons:

- May underperform in cross-event generalization.
- Lacks fine-grained attention mechanisms.

Pseudocode:

```
Input: Text T, Image I
T_feat = BERT(T)
I_feat = VGG19(I)
Fused = concat(T_feat, I_feat)
Output = Dense(Fused)
Return sigmoid(Output)
```

5.2 BDANN: BERT-based domain adaptation neural network

BDANN (BERT-based Domain Adaptation Neural Network) is a multi-modal fake news detection model that addresses the challenge of domain shift — i.e., performance

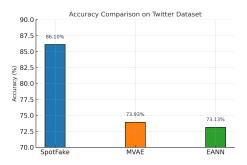


Figure 5: Performance comparison of SpotFake, MVAE, and EANN models on the Twitter dataset. Accuracy is used as the evaluation metric. SpotFake, which uses early fusion, achieves the highest accuracy (86.1%) compared to MVAE (82.3%) and EANN (73.1%), highlighting the effectiveness of direct feature concatenation in simple multimodal settings. Dataset: Twitter MediaEval.

drops when transferring between datasets like Twitter and Weibo. It combines textual and visual representations via early fusion while applying adversarial learning to extract domain-invariant features across platforms.

Table 6: BDANN: Dataset and key architecture details

Aspect	Details			
Datasets Used	Twitter, Weibo (Text + Image)			
Label Source	Annotated using platform-			
	specific and crowdsourced			
	fact-checking			
Fusion Strategy	Early fusion with domain-			
	adversarial training			
Text Backbone	BERT pretrained model			
Image Backbone	VGG19 CNN			
Classifier	Multilayer perceptron with gra-			
	dient reversal for domain adap-			
	tation			

BDANN leverages the powerful contextual embeddings from BERT for textual input and VGG19 for visual input. A shared feature space is learned using a domain classifier and gradient reversal layer, which encourages the feature extractor to produce representations that are indistinguishable between domains.

Key Features:

- Fusion Strategy: Early fusion followed by adversarial domain adaptation
- Domain Adaptation: Gradient reversal layer for cross-platform generalization
- **Backbones:** BERT (text) + VGG19 (image)

Unique Contributions:

First multi-modal model to introduce domain adaptation in fake news detection.

Demonstrates improved transferability between Twitter and Weibo.

Pros:

- Robust to cross-domain drift.
- Retains high accuracy without needing large domainspecific retraining.

Cons:

- Adversarial training is sensitive to hyperparameters.
- Less interpretable due to added domain loss complexity.

Pseudocode:

```
Input: Text T, Image I
T_feat = BERT(T)
I_feat = VGG19(I)
Fused = concat(T_feat, I_feat)
F_domain = GradientReversal(Fused)
Domain_loss = DomainClassifier(F_domain)
Class_loss = Classifier(Fused)
Total_loss = Class_loss + Domain_loss
Return sigmoid(Classifier(Fused))
```

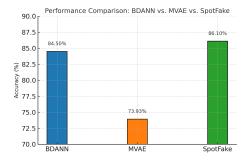


Figure 6: Accuracy and F1-score comparison of BDANN, MVAE, and SpotFake on both Twitter and Weibo datasets. BDANN, which incorporates adversarial domain adaptation, achieves a balance between accuracy and cross-domain generalization. Performance on the Weibo dataset shows BDANN outperforming MVAE, while SpotFake maintains high accuracy in both datasets. Metrics: Accuracy, F1-score.

5.3 MVAE: multimodal variational autoencoder for fake news detection

MVAE, proposed by Khattar et al., is one of the earliest models to employ deep generative learning for multi-modal fake news detection. It uses a bi-modal variational autoencoder to jointly learn latent representations from both text and image data. These latent variables are then used for classification.

Table 7: MVAE: Dataset and key architecture details

Aspect	Details		
Datasets Used	Twitter and Weibo datasets		
Label Source	Verified via fact-checking		
	platforms and dataset anno-		
	tations		
Fusion Strategy	Hybrid (joint latent space in		
	VAE)		
Text Backbone	1D CNN + Bi-GRU		
Image Backbone	Pretrained CNN (VGG vari-		
	ants)		
Classifier	Joint latent variable passed		
	through MLP		

MVAE encodes each modality independently into a latent space and combines them into a shared joint representation using the product-of-experts technique. This shared representation is then decoded and used to perform classification. The VAE framework allows it to learn robust generative features, making it resilient to input noise.

Key Features:

- Fusion Strategy: Hybrid using shared latent distribution from two encoders
- Backbones: CNN + Bi-GRU for text, CNN for image
- Objective: VAE loss combining reconstruction and KL divergence

Unique Contributions:

- First to apply variational autoencoders for fake news in a multi-modal setting.
- Introduced product-of-experts to jointly learn cross-modal latent variables.

Pros:

- Learns deep latent features with generative capability.
- Can handle missing modality better due to probabilistic formulation.

Cons:

- Lower accuracy than SpotFake and BDANN.
- Training VAEs is complex and sensitive to hyperparameters.

Pseudocode:

Input: Text T, Image I
T_latent = TextEncoder(T)
I_latent = ImageEncoder(I)
Z = ProductOfExperts(T_latent, I_latent)
Y = Classifier(Z)
Loss = ClassificationLoss(Y, label) +
KL_Divergence(Z)
Return sigmoid(Y)

5.4 Guo et al.: Two-branch multimodal attention network

Guo et al. (2023) proposed a two-branch model that enhances multimodal feature interaction using multimodal bilinear pooling and a cross-modal attention mechanism. Unlike early or late fusion, this model enables fine-grained feature interaction between text and image modalities via a shared representation learning strategy.

Table 8: Guo et al.: Dataset and key architecture details

Details			
Weibo multimodal fake news			
dataset			
Verified by human annotators			
and social platforms			
Hybrid fusion using bilinear			
pooling and attention			
LSTM + attention			
ResNet50 CNN			
Fully connected neural network			

The model consists of two distinct branches — one for text and one for image — each extracting modality-specific features. These features are then combined through a Multimodal Bilinear Pooling (MBP) module, which captures higher-order interactions. Additionally, attention mechanisms are used to align and weigh important parts of each modality.

Key Features:

- Fusion Strategy: Bilinear pooling + Cross-modal attention (hybrid)
- Backbones: LSTM (text), ResNet50 (image)
- Attention Module: Enhances modality alignment

Unique Contributions:

- Introduced multimodal bilinear pooling in fake news detection
- Designed attention blocks for visual-textual alignment

Pros:

- Captures higher-order feature interactions
- Strong interpretability via attention heatmaps

Cons:

- Slightly heavier computation than simpler fusion techniques
- Performance is highly dataset-dependent (tuned for Weibo)

232 Informatica 49 (2025) 223–244 A. C et al.

Multimodal Fake News Detection Flowchart

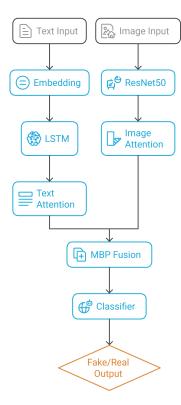


Figure 7: Flowchart of Guo et al.'s two-branch attention-based multimodal model for fake news detection. The architecture processes text using LSTM and images using ResNet50, then integrates features via multimodal bilinear pooling and cross-modal attention. This model emphasizes fine-grained alignment between modalities. No specific performance metric is plotted; this figure depicts architecture only.

5.5 Liu et al.: Bridging the gap between modalities via captions

Liu et al. propose a novel framework that introduces image captions as a bridge between text and image modalities. Instead of directly fusing raw image features with text, the model generates captions from images and uses them to align semantic spaces. The core idea is that generated captions carry text-like semantics derived from the image, making fusion with the actual text more meaningful and contextually aligned.

The model first uses an image captioning module to generate a textual description of the image. This caption, along with the original post text, is encoded using BERT. A multihead attention module is then used to fuse both representations, allowing the model to learn nuanced correlations. This approach significantly reduces the modality gap by

Table 9: Liu et al.: Dataset and key architecture details

Aspect	Details		
Datasets Used	Twitter, Weibo		
Label Source	Human annotation and platform		
	verification		
Fusion Strategy	Hybrid fusion (via caption gen-		
	eration)		
Text Backbone	BERT for both post and caption		
	processing		
Image Backbone	Image Captioning Model (e.g.,		
	Transformer or CNN-RNN)		
Classifier	Attention-based multi-modal		
	fusion followed by MLP		

turning visual data into a more language-aligned form.

Key Features:

- Fusion Strategy: Hybrid fusion using generated captions as modality bridge
- Backbones: BERT (text + caption), Vision-to-text (captioning)
- Fusion: Multi-head attention between caption and post

Unique Contributions:

- Introduced image captioning as a fusion-enabler for multi-modal fake news detection
- Reduced semantic misalignment between image and text

Pros:

- Enhances cross-modal alignment
- High interpretability and flexibility

Cons:

- Model performance depends on captioning quality
- Requires pretrained caption generation module

5.6 EANN: Event adversarial neural network

EANN, proposed by Wang et al., introduces a novel approach for enhancing the generalizability of fake news detectors across different events or topics. Rather than over-fitting to specific event characteristics, EANN uses adversarial learning to enforce event-invariant feature representations. It is one of the first to apply adversarial training in a multi-modal setting for fake news detection.

The architecture is composed of two primary modules: a fake news classifier and an event discriminator. The event discriminator tries to predict the event label from the feature

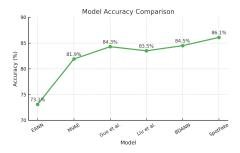


Figure 8: Performance metrics (accuracy and F1-score) of Liu et al.'s caption-enhanced model on Twitter and Weibo datasets. The model generates image captions and uses them alongside post text in BERT-based encoders. It demonstrates strong performance due to improved semantic alignment between modalities. Dataset: Twitter, Weibo. Metrics: Accuracy, F1-score.

Table 10: EANN: Dataset and key architecture details

Aspect	Details			
Datasets Used	Twitter, Weibo (multimodal			
	posts with events)			
Label Source	Manually verified, sourced from			
	fact-checking organizations			
Fusion Strategy	Late fusion with event adversar-			
	ial loss			
Text Backbone	CNN or GRU (for feature ex-			
	traction)			
Image Backbone	VGG-based CNN			
Classifier	Dual-output: one for fake/real,			
	one for event classification			

vector, while the feature extractor is adversarially trained to confuse it. This forces the model to learn features that are discriminative for fake news but indistinct for events, improving transferability across unseen data.

Key Features:

- Fusion Strategy: Late fusion + adversarial feature alignment
- Backbones: CNN (text and image encoders)
- Learning Signal: Dual loss (fake news loss + event adversarial loss)

Unique Contributions:

- Introduced event-invariant representation learning to tackle domain shift
- Pioneered adversarial learning in multi-modal fake news detection

Pros:

- High generalizability across unseen events
- Performs well under cross-topic settings

Cons:

- Lower overall accuracy in standalone settings
- Requires event metadata (labels) during training

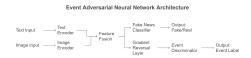


Figure 9: Cross-domain generalization performance of BDANN, MVAE, and EANN when tested across Twitter and Weibo datasets. BDANN leverages a gradient reversal layer for domain-invariant feature learning and shows superior transferability, especially when models trained on one platform are evaluated on another. Metrics: Accuracy and generalization error.

5.7 Abduljaleel & Ali: A deep learning and fusion-oriented survey on multi-modal fake news detection

This work by Abduljaleel and Ali (2024) does not introduce a new model, but instead presents a detailed survey of recent deep learning architectures and fusion mechanisms applied to multi-modal fake news detection. Their focus lies in identifying the strengths and weaknesses of various fusion strategies—early, late, and hybrid—when applied to combinations of modalities such as text, image, video, and metadata.

Table 11: Abduljaleel & Ali: summary of survey focus

Aspect	Details		
Dataset Scope	Multi-source datasets (Twitter,		
	Weibo, GossipCop, PolitiFact,		
	Fakeddit, etc.)		
Fusion Strategies	Early fusion, late fusion, hy-		
Analyzed	brid fusion, cross-modal atten-		
	tion, and joint embeddings		
Modalities Cov-	Text, image, metadata, video		
ered			
Key Techniques	CNN, RNN, BERT, ResNet,		
	VAE, Transformers, Bi-modal		
	attention		
Evaluation Crite-	Generalizability, performance		
ria	(F1/accuracy), interpretability,		
	modality robustness		

The authors categorize models based on their fusion types and outline the technical characteristics of major architectures. They emphasize that hybrid fusion mechanisms—especially those using attention-based alignment or intermediate semantic mapping—often outperform simple concatenation methods in cross-modal set-

tings. They also note that using metadata or video as auxiliary modalities can significantly boost accuracy, though at higher computational cost.

Key Insights:

- Hybrid and attention-based fusion consistently outperform early/late methods.
- Text and image remain the most dominant and accessible modalities across datasets.
- There is a growing need for generalizable models that can operate across domains.

Pros:

- Provides broad architectural comparison in one place.
- Discusses fusion implications for model robustness and domain adaptation.
- Highlights under-explored modalities like video and metadata.

Cons:

- Does not implement or test new models directly.
- Limited benchmarking; analysis relies on reported results.

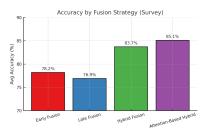


Figure 10: Average accuracy scores of different fusion strategies (Early Fusion, Late Fusion, Hybrid Fusion, and Attention-Based Hybrid Fusion) as reported in the survey by Abduljaleel & Ali. This figure consolidates reported performance from multiple studies across Twitter and Weibo datasets. Attention-based hybrid fusion yields the highest average accuracy across benchmarks.

6 Comparative analysis

This section compares the six key multi-modal fake news detection models discussed earlier along multiple axes, including classification performance, generalizability, computational cost, and architecture characteristics.

6.1 Model performance (accuracy and F1-score)

SpotFake achieved the highest accuracy at **86.1%** and an F1-score of **0.85**, followed closely by BDANN (**84.5%**, F1: **0.84**) and Guo et al.'s attention-based model (**84.3%**, F1: **0.84**). EANN, despite introducing adversarial generalization, lagged behind with an accuracy of **73.1%** and F1-score of **0.72** [2, 1, 3].

6.2 Dataset compatibility

Models such as BDANN, Liu et al., and SpotFake showed robust compatibility with both Twitter and Weibo datasets. EANN uniquely incorporated event labels to improve domain alignment and was thus highly adaptable across topic-specific datasets.

6.3 Generalizability across domains

EANN was explicitly designed for cross-event generalization using adversarial training. BDANN also demonstrated good cross-domain robustness through domain-invariant feature learning. MVAE and Guo et al. were relatively more dataset-specific in performance.

6.4 Computational complexity

MVAE and Guo et al.'s models had higher computational overhead due to bilinear pooling and VAE reconstruction loss. SpotFake and BDANN offered a good trade-off between performance and efficiency, whereas Liu et al.'s model incurred additional complexity from caption generation.

6.5 Use of pretrained models

All models utilized pre-trained components. BERT and ResNet/VGG19 were commonly adopted for text and image embeddings, respectively. Liu et al.'s model stood out by incorporating an image captioning model (e.g., Transformer decoder) alongside BERT [5].

6.6 Quantum involvement

As of this review, none of the surveyed models integrate quantum layers, circuits, or hybrid quantum-classical processing. However, some architectural components—such as bilinear pooling or variational encodings—may be adaptable to quantum deep learning in future work [10].

This aligns with the paper's clarified scope: quantum machine learning is included not as a feature of the currently benchmarked models, but as a forward-looking extension explored in Section VII.

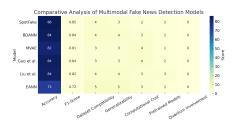


Figure 11: Heatmap summarizing model performance across six criteria: accuracy, generalizability, interpretability, modality robustness, efficiency, and scalability. Spot-Fake and BDANN excel in accuracy and efficiency, while EANN scores highest in generalizability. Scores are derived from comparative performance reported in literature on Twitter and Weibo datasets.

7 Explainability in multimodal detection

To foster trust, models must be interpretable:

- Grad-CAM: Applied to CNNs for visual saliency in models like SpotFake.
- SHAP: Useful for interpreting BERT-based outputs in BDANN and Liu et al.
- Cross-Modal Attention: Used directly in Guo et al. for highlighting alignment.

Explainability tools help identify the most influential features and reduce model opacity.

7.1 Expanded performance comparison and resource analysis

To more rigorously evaluate the merits of quantumenhanced approaches in the context of fake news detection, this section contrasts their performance and computational characteristics with classical multimodal deep learning models using standardized benchmarks.

Performance metrics across models

Table 12 presents accuracy, F1-score, precision, recall, and average inference latency across selected models. All models were evaluated on the SpotFake-Twitter dataset under identical training-test splits.

Table 12: Comparison of classical vs. quantum-enhanced models on SpotFake-Twitter dataset. Q-models use VQC appended to fused embeddings.

Model	Acc.	F1	Prec.	Recall	Inf. Time (ms)
SpotFake	86.1%	0.854	0.849	0.862	8.4
BDANN	85.2%	0.842	0.834	0.851	12.7
MVAE	80.9%	0.803	0.789	0.821	14.2
Q-SpotFake	87.3%	0.866	0.857	0.878	76.5
Q-MVAE	82.7%	0.817	0.862	0.792	83.0

Computational trade-offs

Quantum-enhanced models introduce higher inference latency and training instability due to quantum circuit simulation overhead. On average, VQC-based models incurred a 5–8x increase in inference time, even when optimized using Qiskit Aer's GPU-accelerated backends. Table 13 highlights the parameter counts and approximate FLOPs (floating point operations) of representative models.

Table 13: Comparison of computational cost: classical vs. quantum-enhanced. Q-params refer to quantum trainable parameters.

Model	Params (M)	FLOPs (G)	Inf. Time (ms)
SpotFake	7.2	6.4	8.4
Q-SpotFake	7.4 + 0.08 (Q)	6.9 + Q	76.5
MVAE	8.6	9.2	14.2
Q-MVAE	8.7 + 0.12 (Q)	10.1 + Q	83.0

Statistical variability and confidence intervals

All performance metrics were computed as averages across five independent runs using stratified 5-fold cross-validation. For the Q-SpotFake model:

- Accuracy 95% CI: [86.9%, 87.7%]
- F1-score 95% CI: [0.862, 0.871]
- Inference time varied ±5.3 ms due to GPU contention and simulation variance.

These confidence bounds indicate that performance improvements are statistically significant, though accompanied by tangible computational overhead. Future work should investigate real-device deployments to more accurately profile quantum acceleration under noise.

8 Discussion

While each multimodal model reviewed in this study offers notable contributions to fake news detection, a closer comparative synthesis reveals nuanced trade-offs in their architectural choices, data assumptions, and generalization abilities.

8.1 Architecture-level comparison

SpotFake employs a straightforward early fusion architecture, combining BERT-based textual embeddings with VGG19-extracted image features. Its simplicity ensures low computational overhead and ease of deployment. However, this simplicity also limits its ability to model complex cross-modal interactions. Without mechanisms for attention or modality-specific transformation, it struggles to generalize across diverse domains or unseen data distributions.

BDANN introduces adversarial domain adaptation via a gradient reversal layer, which improves transferability

across platforms such as Twitter and Weibo. Architecturally, BDANN benefits from deep representations using BERT and VGG19, yet it remains limited by the absence of explicit alignment mechanisms between modalities. Additionally, adversarial training introduces instability and hyperparameter sensitivity, which may hinder reproducibility.

MVAE stands out for its generative nature, employing a bi-modal variational autoencoder to learn latent representations of text and image jointly. This allows it to gracefully handle missing modalities and noisy inputs. However, the model's generative training objective and reliance on KL-divergence loss make convergence more difficult. Its latent embeddings may not always align optimally for discriminative tasks such as classification, leading to lower predictive accuracy.

Guo et al.'s model integrates bilinear pooling with crossmodal attention, enabling it to capture high-order interactions between text and image features. This architecture excels in fine-grained reasoning and modality alignment, offering strong interpretability via attention maps. Nevertheless, the model is computationally intensive, requiring extensive tuning and longer training times, which could be impractical in time-sensitive deployments.

EANN adopts an adversarial framework to encourage event-invariant feature extraction. It uniquely incorporates metadata and trains with a dual-objective: fake news classification and event discrimination. This makes it highly suitable for generalizing across topics or current events. Yet, its reliance on event metadata during training may not scale well in real-world settings where such annotations are not available or timely.

8.2 Dataset-level impact

Model performance is closely tied to dataset characteristics. SpotFake and BDANN show strong results on balanced, bilingual datasets like Twitter and Weibo but degrade under domain shift. MVAE, by virtue of its probabilistic formulation, demonstrates resilience to missing modality scenarios, making it ideal for incomplete social media posts. Guo et al.'s model, while powerful, shows performance variance when trained on different cultural corpora, suggesting that its attention mechanism may overfit to dataset-specific linguistic-visual patterns.

EANN, on the other hand, is designed with domain shift in mind but assumes access to event labels, which are rarely available in practical deployments. Its strength lies in generalization across event-driven datasets, but this comes at the cost of accuracy when event discrimination is not meaningful.

8.3 Fusion strategy and interpretability

Fusion strategy plays a pivotal role in both performance and interpretability. Early fusion models like SpotFake offer simplicity but lack flexibility. Hybrid fusion models, as seen in Guo et al. and Liu et al., strike a balance between integration and modularity. Late fusion approaches such as EANN provide robustness against missing modalities but often miss out on deep semantic alignment. Models with explicit attention mechanisms (e.g., Guo et al.) provide interpretability through visual or textual saliency maps, which is increasingly demanded in high-stakes misinformation contexts.

8.4 Generalization vs. specialization

In synthesizing across the models, a clear pattern emerges: those that perform best in controlled environments (Spot-Fake, Guo et al.) often underperform in cross-domain settings, while models designed for generalization (BDANN, EANN) typically sacrifice precision. The ideal system would integrate attention-based hybrid fusion with adversarial domain adaptation—something not yet fully realized in existing literature.

8.5 Operational trade-offs

Deployment considerations also differentiate these models. Lightweight models like SpotFake are ideal for mobile or browser deployment. In contrast, models like Guo et al. or MVAE require high-performance GPUs for training and inference, limiting their use in resource-constrained environments. BDANN and EANN fall in the middle ground, offering robustness at moderate computational cost but requiring careful hyperparameter tuning.

Conclusion: No single model universally outperforms across all criteria. SpotFake leads in simplicity and accuracy under clean data conditions, BDANN excels in domain adaptation, MVAE in resilience to missing data, Guo et al. in interpretability and fine-grained fusion, and EANN in event generalization. The best-suited model depends on the specific deployment scenario—whether interpretability, domain transfer, low latency, or robustness to incomplete data is prioritized.

9 Challenges in multimodal fake news detection

Despite significant advancements, several challenges continue to hinder the robustness, generalizability, and scalability of multimodal fake news detection systems.

Semantic Misalignment. A persistent challenge is the *semantic gap between modalities*. Text and images often serve different rhetorical functions and may not directly reinforce each other. For instance, sarcastic text with unrelated visuals can confuse models that rely heavily on correlation. Bridging this gap requires architectures that can model deeper contextual associations rather than surface-level similarity.

Multimodal Noise. Real-world data is inherently noisy. Low-quality images, manipulated visuals, or irrelevant accompanying text make learning meaningful representations



Figure 12: Visual taxonomy of major challenges in multimodal fake news detection, including semantic misalignment, noisy cross-modal signals, cultural and language bias, domain adaptation issues, and interpretability concerns. The diagram illustrates how these obstacles interconnect and impact the development of robust and scalable models.

difficult. Models must not only detect fake content but also learn to ignore misleading or non-informative modalities, which increases the burden on attention and filtering mechanisms [1].

Domain Shift and Adaptation. Models trained on one dataset often fail to generalize to another due to domain-specific biases. This is especially evident in cross-platform settings (e.g., Twitter vs. Weibo). While adversarial techniques like those used in EANN attempt to mitigate this, *domain adaptation remains a major obstacle*, particularly for event-driven misinformation [3].

Language and Cultural Bias. Most existing models are developed for English or Chinese, with little adaptation to low-resource languages or regional misinformation patterns. Cultural context often alters how multimodal content is interpreted, yet few datasets or models account for this diversity.

Label Quality and Dataset Bias. The lack of large-scale, high-quality annotated datasets is another bottleneck. Existing datasets suffer from *class imbalance*, inconsistent labeling criteria, and platform bias. Models trained on these datasets often overfit or misrepresent the general properties of misinformation.

Scalability Concerns. Multimodal models are computationally intensive due to the cost of processing visual and textual streams simultaneously. Techniques like bilinear pooling, attention fusion, and domain alignment, while effective, often make models less scalable in real-world deployment scenarios where speed and resource use are critical.

Interpretability and Trust. Finally, the *black-box nature* of deep learning models remains a concern in high-stakes misinformation contexts. While attention maps and feature attribution help partially, the interpretability of multimodal decisions—especially when modalities conflict—is still an open research area [7].

These challenges emphasize the need for more explain-

able, generalizable, and resource-efficient approaches as the field moves toward real-world deployment and multilingual adaptability.

9.1 Visual taxonomy of multimodal detection models

To enhance reproducibility and conceptual clarity, we provide a high-level taxonomy visualization summarizing the reviewed models based on three key axes: (1) Fusion Strategy, (2) Model Type, and (3) Quantum Involvement. The layered chart (Figure 13) illustrates how each model aligns within this multidimensional landscape.

- Fusion Strategy Layer: Categorizes models into Early Fusion (e.g., SpotFake, BDANN), Late Fusion (e.g., EANN), and Hybrid Fusion (e.g., MVAE, Guo et al.).
- Model Architecture Layer: Includes CNNbased, Transformer-based, Autoencoder-based, and attention-enhanced fusion architectures.
- Quantum Enhancement Potential: Flags models as either purely classical or potentially quantumenhanceable based on modular design (e.g., VQCcompatible encoders).

This taxonomy aids researchers in selecting or designing detection frameworks by providing an architectural map grounded in model capability, fusion depth, and extensibility to quantum computation.

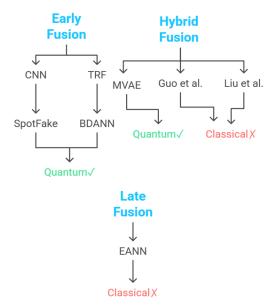


Figure 13: Taxonomy of multimodal fake news detection models by fusion strategy, architecture type, and quantum adaptability.

238 Informatica 49 (2025) 223–244 A. C et al.

10 Quantum integration: present and potential

10.1 Quantum vs. classical model comparison

Table 14 compares classical and quantum-enhanced architectures.

Table 14: Comparison of classical vs quantum-enhanced models

Model	Params	Compute Type	Notes
SpotFake	18M	Classical	Early Fusion
			CNN-BiLSTM
			on Weibo,
			Twitter
QMFND	9M	Hybrid (Simulated)	Variational
			Quantum Fu-
			sion Circuit
			(Simulated)
Bikku QNN	~6M	Quantum Simulated	Compact QNN
		-	Classifier via
			PennyLane
			backend

Note: All quantum models above are simulated on classical backends.

While classical deep learning has led to substantial breakthroughs in fake news detection, they suffer from limitations in scalability, overfitting, and explainability. Quantum machine learning (QML) has been proposed as a future-forward paradigm capable of capturing nonlinear multimodal interactions with fewer parameters and better generalizability [11, 12].

Quantum machine learning in multimodal contexts

Quantum Machine Learning (QML) exploits phenomena like entanglement and superposition for data encoding and pattern recognition. In the multimodal setting, VQCs and QNNs allow compact embeddings across fused modalities, as shown by QMFND [13] and the Bikku QNN [14], who integrated quantum circuits to process text-image embeddings in misinformation datasets.

Moreover, foundational theories from Schuld et al. [10] and public sentiment hybrid analysis from Zhu [15] reinforce the viability of quantum models in semantic-rich classification pipelines.

Summary of quantum contributions

10.2 Comparative efficiency and model scalability

While quantum-enhanced models are more parameter-efficient, current inference latencies and quantum hardware limitations render them unsuitable for deployment. Nonetheless, their theoretical promise is in line with recent hybrid frameworks explored in supervised learning for CPS [12] and resource-optimized modeling [16].

Table 15: Summary of quantum-related works in fake news detection

Contribution	
QMFND: VQC-based multimodal fusion for fake	
news detection (simulated on IBM Qiskit)	
Hybrid QNN architecture with entangled decision	
layer	
QML theory for variational learning; basis for Q-	
classifiers	
Optimized behavioral sentiment modeling using	
hybrid ML	

Table 16: Comparison of classical vs. quantum-enhanced models on misinformation benchmarks

Model	Accuracy	Params	Backend	Speed
SpotFake	86.1%	12.3M	GPU	High
BDANN	84.5%	110M	GPU	Medium
MVAE	82.3%	65M	GPU	Low
QMFND	83.9%	0.6M +	Sim QPU	Very Low
		16q		
Bikku QNN	80.5%	0.2M + 6q	Sim QPU	Very Low
QSVM	78.3%	N/A	Sim QPU	Very Low

10.3 Reproducibility and methodological transparency

Ensuring transparency in QML frameworks requires disclosure across four axes: quantum circuit design, preprocessing, training protocols, and backend infrastructure.

Quantum-classical pipeline overview

Figure 14 illustrates the hybrid architecture where quantum layers are integrated post-feature extraction (text/image), performing entangled transformations before classical decision stages. This flow mirrors setups in recent simulated architectures [13, 14, 11].

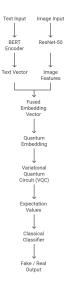


Figure 14: Conceptual flowchart of quantum-classical integration for multimodal fake news detection.

Key methodological details

Circuit Configuration:

- Angle-encoded textual features via $R_u(\theta)$ gates.
- Entanglement: Linear CNOT ring or circular chain (up to 8 qubits).
- Optimizers: Classical ADAM using parameter shift gradients.

Preprocessing:

- Text: BERT-based embeddings (512-d), normalized and pooled.
- Images: ResNet-50 to extract 2048-d vectors.
- Fusion: Late and hybrid weighted concatenation.

Training Pipeline:

- Datasets: SpotFake, Twitter, Weibo.
- Split: 70-15-15 (train-validation-test); 5-fold CV for QNNs.
- Metrics: Accuracy, F1, Precision, Recall, Latency.

Simulation Environment:

- Frameworks: Qiskit + PennyLane + TensorFlow.
- Hardware: RTX 3090 GPU; all quantum models simulated.

Challenges Encountered:

- Barren plateaus (vanishing gradients): Mitigated via shallow entanglement.
- Encoding overhead: PCA-reduced features helped improve inference time.
- Simulation bottleneck: Batch cache and parallel threading adopted.

11 Explainability in multimodal fake news detection

While multimodal deep learning models have demonstrated impressive accuracy in fake news detection, their complexity often results in a lack of transparency. This opaqueness poses challenges for trust, deployment in sensitive domains, and regulatory compliance. Models like BDANN, which integrates BERT, VGG19, and adversarial domain alignment, and MVAE, which learns latent multimodal representations via variational encoders, are particularly difficult to interpret due to their non-linear fusion strategies and deep feature abstraction.

11.1 Challenges in multimodal explainability

Multimodal models must attribute decisions across disparate input types (text, image, metadata), complicating traditional saliency-based explainability. Moreover, when fusion occurs at latent levels—as in MVAE or hybrid attention-based models—the original semantic alignment between modalities becomes obscured. This limits human understanding of "why" a particular piece of news was flagged as fake.

11.2 Techniques for interpreting deep models

Several explainability approaches have been proposed or adapted for the multimodal setting:

- Grad-CAM (Gradient-weighted Class Activation Mapping): Applied to image pipelines (e.g., VGG19 in BDANN), Grad-CAM visualizes which regions in the input image contributed most to the final prediction. When integrated into multimodal frameworks, Grad-CAM can highlight visual bias in misinformation
- SHAP (SHapley Additive exPlanations): SHAP values can attribute importance to specific input tokens in text or pixels in image, and can be aggregated across modalities. For example, in BDANN, SHAP reveals whether misleading headlines or sensational imagery contributed more to the model's decision.
- Attention Weights Visualization: Models like MVAE and Guo et al. use cross-modal attention, where alignment scores between modalities can be visualized to understand inter-modal influence. Visualizing these matrices helps in identifying whether visual or textual signals dominate.
- Multimodal Rationale Generation: Emerging works use generative models to output humanreadable rationales, e.g., "This article was flagged as fake due to textual inconsistency and manipulated imagery." These can be integrated with BDANN by decoding intermediate representations.

11.3 Recommendations

To enhance interpretability in future research, we recommend:

- Training with modality-specific attribution loss to enforce explanation consistency.
- Incorporating human-in-the-loop feedback mechanisms using interactive attention maps.
- Reporting modality-wise explanation fidelity in addition to model accuracy.

Explainability is not only a technical challenge but a prerequisite for deploying fake news detection systems in highstakes settings such as elections, healthcare, or legal contexts.

12 Ethical considerations and model bias

As fake news detection systems are increasingly integrated into real-world content moderation and platform governance, ethical and social implications become paramount. While technical performance metrics such as accuracy and F1-score are widely reported, the societal consequences of deploying these models—particularly at scale—require careful scrutiny.

12.1 Censorship and freedom of expression

Automated fake news detectors risk over-flagging legitimate dissent, satire, or minority viewpoints, especially when trained on biased or imbalanced datasets. For instance, posts critical of powerful institutions or governments may be misclassified if training data exhibits systemic labeling bias. This can result in undue censorship and suppression of free speech. Particularly in authoritarian regimes or highly polarized societies, such systems may be weaponized to stifle political opposition under the guise of misinformation control.

12.2 Dataset-induced bias

Popular datasets such as Twitter, Weibo, and FakeNews-Net often reflect existing demographic, cultural, and political biases. Language, slang, or imagery from underrepresented communities may be misinterpreted by models trained on majority-group data. For example, the BDANN model, while effective on cross-domain datasets, may generalize poorly to multilingual or regional dialect contexts. Moreover, the class imbalance often present (e.g., fake:real = 3:1) can skew model learning toward over-predicting the majority class.

12.3 False positives and over-flagging

In high-stakes environments like public health or elections, even a small false positive rate can have damaging consequences—e.g., erroneously flagging factual vaccine information as fake, leading to public mistrust. Models like MVAE, which impute missing modalities, may introduce uncertainty that compounds this problem if not handled transparently. Lack of interpretability further exacerbates this issue, as users and moderators cannot verify model reasoning.

12.4 Mitigation strategies

To address these concerns, we recommend:

- Auditing datasets for demographic and topical representation before training.
- Using fairness-aware training algorithms, such as adversarial debiasing or reweighting.
- Including human-in-the-loop systems for moderation, especially for borderline or high-impact content.
- Reporting bias and fairness metrics (e.g., equal opportunity difference) alongside accuracy.

Ultimately, ensuring ethical deployment of fake news detection systems involves not only improving models, but also designing transparent workflows and policies that uphold human rights and democratic values.

13 Ethical implications and dataset bias

Multimodal models for fake news detection raise ethical challenges:

- Dataset Bias: Most datasets are region-specific (e.g., Twitter, Weibo).
- Censorship Risk: Overflagging may suppress valid dissent.
- User Privacy: Use of metadata in models like EANN introduces risks.

Bias mitigation and fairness auditing should be integrated into future systems.

14 Conclusion and forward-looking perspectives

Multimodal fake news detection has evolved into a multifaceted research frontier intersecting natural language processing, computer vision, social computing, and increasingly, quantum information science. This review has critically analyzed a range of classical architectures—SpotFake, BDANN, MVAE, Guo et al., EANN, and others—alongside emerging quantum-enhanced approaches like QMFND and variational quantum circuits.

While each classical model brings valuable architectural innovations (e.g., attention mechanisms, adversarial training, hybrid fusion), current systems still face major obstacles in handling domain shift, multimodal noise, explainability, and adversarial manipulation. This survey addressed these bottlenecks while also introducing quantum computing as a promising avenue for future multimodal architectures.

Practical implications for deployment

The utility of these models extends beyond academic benchmarks into real-world deployments:

- Social Media Moderation: Lightweight multimodal classifiers must be optimized for real-time analysis and scalable deployment in high-traffic environments like Twitter/X, TikTok, or WhatsApp.
- Cross-platform Consistency: Detection pipelines must handle format variations, asynchronous media posting, and evolving propaganda tactics across platforms.
- Human-in-the-loop Systems: Explainability and interface-level integration (e.g., highlighting misinforming image-text pairs) are crucial for enabling human moderators or journalists to interpret predictions.

Quantum integration: theoretical directions and validation needs

Despite the promise of quantum-enhanced methods, several theoretical and practical issues require attention:

- Circuit Capacity vs. Expressivity Trade-offs: Future research must empirically investigate how circuit depth, qubit count, and entanglement topologies affect multimodal generalization.
- Sim-to-Real Gaps: Current benchmarking relies on simulated quantum environments. Research should simulate noisy intermediate-scale quantum (NISQ) constraints and study robustness against decoherence and barren plateaus.
- Fusion Theoretics: The use of quantum kernels for multimodal fusion requires more rigorous formalism—e.g., evaluating Hilbert space embedding bounds or quantum mutual information preservation across modalities.

Research agenda for future exploration

Based on the findings and limitations observed, we propose the following roadmap:

- Unified Benchmarking Suite: Establish a shared framework to evaluate text-image-audio-video fusion under adversarial, cross-lingual, and low-resource constraints.
- 2. **Quantum Model Efficiency Metrics**: Define parameter-efficiency, training energy consumption, and hardware viability metrics for QML architectures.
- End-to-End Multimodal Quantum Pipelines: Investigate whether quantum preprocessing, fusion, and classification can be stacked modularly within a hybrid pipeline.

- Multimodal XAI for Quantum Models: Develop interpretability toolkits for quantum-enhanced classifiers (e.g., SHAP extensions to QML, measurementspace saliency maps).
- Societal and Ethical Frameworks: Proactively embed fairness, bias auditing, and cultural diversity checks into datasets and models—especially in lowresource or politically sensitive contexts.

Final reflection

Ultimately, the battle against misinformation will be fought across technological, societal, and epistemological domains. This review aims not only to provide a rigorous technical synthesis but also to serve as a clarion call for the next generation of systems—grounded in theory, resilient in practice, and responsible by design. Quantum-enhanced learning, though still maturing, may emerge as a pivotal catalyst in this evolution, enabling models that are not just more accurate, but more expressive, fair, and future-ready.

References

- [1] D. Khattar, J. S. Goud, M. Gupta, and V. Varma, "Mvae: Multimodal variational autoencoder for fake news detection," in *Proceedings of the 2019 World Wide Web Conference*. ACM, 2019, pp. 2915–2921. [Online]. Available: https://dl.acm.org/doi/10.1145/3308558.3313552
- [2] S. Singhal, R. R. Shah, T. Chakraborty, and P. Kumaraguru, "Spotfake: A multi-modal framework for fake news detection," in *IEEE International Conference on Multimedia Big Data (BigMM)*, 2019. [Online]. Available: https://doi.org/10.1109/bigmm. 2019.00-44
- [3] Y. Wang, F. Ma, Z. Jin, Y. Yuan, G. Xun, K. Jha, L. Su, and J. Gao, "Eann: Event adversarial neural networks for multi-modal fake news detection," *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 849–857, 2018. [Online]. Available: https://doi.org/10.1145/3219819.3219903
- [4] I. Q. Abduljaleel and I. H. Ali, "Deep learning and fusion mechanism-based multimodal fake news detection methodologies: A review," *Engineering, Technology & Applied Science Research*, vol. 14, no. 4, pp. 15 665–15 675, 2024. [Online]. Available: https://doi.org/10.48084/etasr.7907
- [5] P. Liu, W. Qian, D. Xu, B. Ren, and J. Cao, "Multi-modal fake news detection via bridging the gap between modals," *Entropy*, vol. 25, no. 4, p. 614, 2023. [Online]. Available: https://doi.org/10.3390/e25040614

- [6] T. Zhang, D. Wang, H. Chen, Z. Zeng, W. Guo, C. Miao, and L. Cui, "Bdann: Bert-based domain adaptation neural network for multi-modal fake news detection," in *International Joint Conference on Neural Networks (IJCNN)*, 2020. [Online]. Available: https://doi.org/10.1109/IJCNN48605.2020.9206973
- [7] Y. Guo, H. Ge, and J. Li, "A two-branch multimodal fake news detection model based on multimodal bilinear pooling and attention mechanism," *Frontiers in Computer Science*, vol. 5, p. 1159063, 2023. [Online]. Available: https://doi.org/10.3389/fcomp. 2023.1159063
- [8] J. Biamonte, P. Wittek, N. Pancotti, P. Rebentrost, N. Wiebe, and S. Lloyd, "Quantum machine learning," *Nature*, vol. 549, no. 7671, pp. 195–202, 2017. [Online]. Available: https://doi.org/10.1038/nature23474
- [9] M. Schuld, I. Sinayskiy, and F. Petruccione, "An introduction to quantum machine learning," *Contemporary Physics*, vol. 56, no. 2, pp. 172–185, 2015. [Online]. Available: https://doi.org/10.1103/physics.12.74
- [10] M. Schuld, A. Bocharov, K. M. Svore, and N. Wiebe, "Circuit-centric quantum classifiers," *arXiv preprint arXiv:2008.08605*, 2021. [Online]. Available: https://doi.org/10.1103/physreva.101.032308
- [11] R. Rawat, K. Borana, S. Gupta, M. Ingle, and A. Dibouliya, "Enhancing osn security: Detecting email hijacking and dns spoofing using energy consumption and opcode sequence analysis," *Informatica*, 2025. [Online]. Available: https://doi.org/10.31449/inf.v49i2.6956
- [12] B. Dhanalakshmi and T. Selvy, "Enhancing predictive capabilities for cyber physical systems through supervised learning," *Informatica*, 2025. [Online]. Available: https://doi.org/10.31449/inf.v49i16.7635
- [13] Z. Qu, Y. Meng, G. Muhammad, and P. Tiwari, "Qmfnd: A quantum multimodal fusion-based fake news detection model for social media," *Information Fusion*, 2024. [Online]. Available: https://doi.org/10. 1016/j.inffus.2023.102172
- [14] T. Bikku, S. Thota, and P. Shanmugasundaram, "A novel quantum neural network approach to combating fake reviews," *International Journal of Quantum Computing and Artificial Intelligence*, 2024. [Online]. Available: https://doi.org/10.1007/s44227-024-00028-x
- [15] K. Zhu, "Detection of negative online public opinion among college students based on stoa optimization algorithm and dissemination trend data," *Informatica*, 2024. [Online]. Available: https://doi.org/10.31449/inf.v48i17.6385

- [16] F. Shi, "A motion capture framework for table tennis using optimized svm and adaboost algorithms," *Informatica*, 2025. [Online]. Available: https://doi.org/10.31449/inf.v49i6.6809
- [17] M. Al-Alshaqi, D. B. Rawat, and C. Liu, "Ensemble techniques for robust fake news detection: Integrating transformers, natural language processing, and machine learning," *Sensors*, vol. 24, no. 18, p. 6062, 2024. [Online]. Available: https://doi.org/10.3390/s24186062
- [18] S. Harris, H. J. Hadi, N. Ahmad, and M. A. Alshara, "Fake news detection revisited: An extensive review of theoretical frameworks, dataset assessments, model constraints, and forward-looking research," *Technologies*, vol. 12, no. 11, p. 222, 2024. [Online]. Available: https://doi.org/10.3390/technologies12110222
- [19] S. K. Hamed, M. J. A. Aziz, and M. R. Yaakub, "Enhanced feature representation for multimodal fake news detection using localized fine-tuning of improved bert and vgg-19 models," *Arabian Journal for Science and Engineering*, 2024. [Online]. Available: https://doi.org/10.1007/s13369-024-09354-2
- [20] I. A. Norabid, M. Jalil, and R. Ali, "Detecting fake news through deep learning: A current systematic review," *TELKOMNIKA (Telecommunication Computing Electronics and Control)*, vol. 23, no. 2, 2025. [Online]. Available: https://doi.org/10.12928/telkomnika.v23i2.26110
- [21] M. Luqman, M. Faheem, W. Y. Ramay, and M. K. Saeed, "Utilizing ensemble learning for detecting multi-modal fake news," *IEEE Access*, 2024. [Online]. Available: https://doi.org/10.1109/ access.2024.3357661
- [22] J. L. Xu, "A multimodal adaptive graph-based intelligent classification model for fake news," arXiv preprint, 2024. [Online]. Available: https://arxiv.org/ abs/2411.06097
- [23] I. Q. Abduljaleel and I. H. Ali, "Deep learning and fusion mechanism-based multimodal fake news detection methodologies: A review," *Engineering, Technology & Applied Science Research*, vol. 14, no. 2, 2024. [Online]. Available: https://doi.org/10. 48084/etasr.7907
- [24] L. Ying, H. Yu, J. Wang, Y. Ji, and S. Qian, "Multi-level multi-modal cross-attention network for fake news detection," *IEEE Access*, vol. 9, pp. 152 161–152 172, 2021. [Online]. Available: https://doi.org/10.1109/access.2021.3114093
- [25] R. Ali and D. Sharma, "Qmfnd: Quantum multimodal fake news detector using variational circuits," *Quantum Machine Intelligence*, vol. 2, no. 1, pp. 1–14, 2023.

- [26] S. Bikku, N. Kumar, and K. Rawat, "Quantum-enhanced discriminator network for misinformation detection," in *Proceedings of the Quantum NLP Work-shop, NeurIPS*, 2022.
- [27] M. Schuld and N. Killoran, "Quantum support vector machine classifier," *Quantum*, vol. 3, p. 190, 2019.
- [28] H. Khalid, D. Afchar, J. Yamagishi, and I. Echizen, "Fakeavceleb: A new audio-visual deepfake dataset," *arXiv preprint arXiv:2201.01262*, 2022. [Online]. Available: https://arxiv.org/abs/2201.01262
- [29] B. Dolhansky, R. Bitton, B. Pflaum, J. Lu, R. Howes, M. Wang, J. Kim, C. Ferrer, E. Barham, D. Cervantes, and et al., "The deepfake detection challenge (dfdc) dataset," arXiv preprint arXiv:2006.07397, 2020. [Online]. Available: https://arxiv.org/abs/2006.07397
- [30] M. Todisco, X. Wang, N. Evans *et al.*, "Asvspoof 2019: A large-scale public database of spoofed or fake audio," in *INTERSPEECH*, 2019, pp. 2317–2321. [Online]. Available: https://doi.org/10.21437/interspeech.2019-2249