



ARTICLE

## Causal Cross-Modal Context Fusion for Real-Time Video Summarization with Predictive Tracking and Validated Adaptive Evaluation

Aravapalli Rama Satish<sup>1</sup>, Sai Babu Veeram<sup>2,\*</sup>, Shonak Bansal<sup>3,\*</sup>, Krishna Prakash<sup>4</sup> and Mohammad Rashed Iqbal Faruque<sup>5,\*</sup>

<sup>1</sup>School of Computer Science and Engineering, VIT-AP University, Amaravati, India

<sup>2</sup>Department of AI & DS, KLEF Deemed to be University, Vaddeswaram, Guntur, India

<sup>3</sup>Department of Electronics and Communication Engineering, Chandigarh University, Gharuan, Mohali, India

<sup>4</sup>Department of Information and Communication Technology, Marwadi University, Rajkot, India

<sup>5</sup>Space Science Centre (ANGKASA), Institute of Climate Change (IPI), Universiti Kebangsaan Malaysia, Bangi, Malaysia

\*Corresponding Authors: Sai Babu Veeram. Email: saibabuv@gmail.com; Shonak Bansal. Email: shonakk@gmail.com; Mohammad Rashed Iqbal Faruque. Email: rashed@ukm.edu.my

Received: 09 February 2026; Accepted: 14 May 2026; Published: 30 June 2026

**ABSTRACT:** Real-time video streams now flood everything from security cameras to social media, yet current summarization systems still stumble when audio, visual, and semantic cues unfold with tangled cause-and-effect patterns. Most cross-modal transformers treat correlations as if time were a flat canvas, ignoring how an early sound might trigger a later visual event in the process. They also lack mechanisms to predict tracking uncertainty, adapt to narrative shifts, or evolve their own evaluation criteria, leaving summaries brittle and often incoherent in process. To address these gaps, we propose a Cross-Modal Context Fusion framework built from five tightly linked components. A Temporal-Causal Graph Memory Network captures directional cause-and-effect edges across audio, video, and semantic signals, improving the logical flow of detected key segments. Additionally, a Predictive Entropy Reinforcement Engine learns camera focus and keyframe decisions that minimize future uncertainty, thereby stabilizing tracking under rapid motion or noise. The Cross-Modal Residual Synergy Transformer explicitly models discrepancies, such as off-screen speech, and feeds those residuals back to refine the fusion. For long-form narrative coherence, a Dynamic Hierarchical Context Predictor alternates between micro-actions and macro-story arcs, balancing fine detail with global structure. Finally, a Self-Evolving Evaluation Loop meta-learns to adjust loss weights as deployment contexts shift, sustaining performance without costly full retraining sets. Experiments on SumMe, TVSum, and long-form documentaries indicate up to 15% F1 and 12% ROUGE-L gains, with human studies reporting 18% higher perceived coherence and >90% sustained approval in process. The result is a video summarizer that reasons causally, anticipates uncertainty, adapts its own metrics, and delivers concise yet narratively faithful summaries suited for demanding real-time applications.

**KEYWORDS:** Adaptive evaluation; cross-modal fusion; process; reinforcement learning; temporal-causal graphs; video summarization

### 1 Introduction

The expansion of real-time video streams from surveillance cameras, live broadcasts, autonomous vehicles, and social media platforms has increased the need for concise, semantically rich summaries that can be viewed and acted upon quickly. Keyframe extraction, shot boundary identification, and transformer-style

cross-modal attention struggle with temporal causality and static assessment metrics. Most methods consider visual, semantic, auditory, and cue alignment as a correlation task, neglecting how early sounds may cause later visual events or object motion affect future saliency sets. The summaries often lack logical narrative flow and fail in loud or fast-changing environments. This work proposes a Cross-Modal Context Fusion architecture that strongly integrates causal reasoning, predictive monitoring, and self-evolving evaluation to close this gap in the process. The TCGMN is essential to this design because it records directional cause-and-effect links between auditory, visual, and semantic data samples. Unlike traditional attention approaches, TCGMN uses a graph with delay-aware edges to infer how earlier auditory or visual stimuli affect video saliency sets. Over this causal backbone, a PERE guides camera focus and keyframe selection to minimize object location and audio event prediction uncertainty, maximizing entropy reduction.

To prevent missing subtle but semantically relevant events, the CRST separates and re-injects modalities, including noisy off-screen voices and obscured objects. The DHCP provides narrative summaries using micro-level action modeling and macro-level subject segmentations. Finally, a SEEL meta-learns from user feedback and deployment circumstances to adjust loss weighting and optimization objectives to preserve performance when data distributions change. SumMe, TVSum, and long-form documentary datasets reveal that the proposed system increases F1, ROUGE-L, and human-rated coherence by 15% and 12%, respectively. These findings show that causal cross-modal reasoning, predictive reinforcement, and adaptive evaluation can improve real-time video summary logic and context.

### ***Motivation & Contribution***

Awareness that current video summarizing models are correlation-driven and fragile drove our research. Traditional cross-modal transformers synchronize audio and visual input but rarely apply directed cause-and-effect logic. In surveillance under changing lighting, mobile video streaming with background noise, and sports broadcasting with off-screen aural cues, such approaches fracture narratives and hinder object motions. Fixed assessment measures hinder traditional systems from adapting to changing user expectations or operational constraints. The architecture must incorporate multimodal input, reason about temporal causation, foresee ambiguity, and alter its goals to address these challenges.

This study addresses those needs using a causally grounded cross-modal paradigm and five synergistic components. Time-lagged acoustic, visual, and semantic cause-and-effect relationships are stored in the Temporal-Causal Graph Memory Network. To eliminate uncertainty, the Predictive Entropy Reinforcement Engine rewards future tracking stabilization. Undetected speakers and movements are detected using the Cross-Modal Residual Synergy Transformer. Dynamic Hierarchical Context Predictor models micro actions and macro narrative arcs to make summaries more than salient shots. The Self-Evolving Evaluation Loop uses changing conditions to meta-learn optimal metric weightings to close the adaptation gap. These modules enhance F1 and ROUGE-L scores, tracking drift and human-perceived coherence while laying the groundwork for next-generation real-time video summarization systems that learn, reason, and adapt in the process.

## **2 Review of Existing Models Used for Video-Based Human Feature Analysis**

From heuristic methods to deep learning and cross-modal systems that emphasize causality, narrative flow, and real-time adaptability, video summarizing research has evolved rapidly in recent years. A review of important contributions analysis reveals this evolution and repeated themes that drive the current state of the art. A broad conceptual framework is followed by federated learning, reinforcement-driven tracking, multimodal fusion, and health, sports, and security-specific summarizations. An overview by Kadam and Deshpande [1], an important foundation is cataloging query-attentive video summarizing techniques and showing how early systems relied on attention processes but lacked comprehensive temporal causality

management [1]. Based on this basis, Aravinda et al. [2] forecast frames using a convolutional LSTM and 3D CNN, a temporal continuity-modeling architecture change. Kandaswamy and Balachandern [3] study a self-gated federated capsule network for multi-view summarization and privacy-preserving distributed training. Khalid et al. [4] use fuzzy C-mean clustering for visual feature fusion, while Yarrarapu et al. [5] focus on items of interest using MobileNetSSD. Ravishankar et al. [6] focus on compression and saliency-aware spatial-temporal integration process. Babu Veeram and Satish [7] provide a statistical foundation that supports the development of the Causal Cross-Modal Context Fusion.

Iteratively, Next, as per Table 1, Deepa et al. [8] use graph neural networks for dynamic summarizing, while Guan et al. [9] discuss multi-object tracking, a foundation for many summarization pipelines. Some early studies enhance spatiotemporal graph understanding. Recursive spatiotemporal graph modeling with language direction by Park et al. [10] provides natural language summaries. Evolutionary algorithms and deep learning to summarize static films. Öztürk et al. use linguistic summarizing to analyze autism data in healthcare [11]. Coauthor Blanco-Fernández et al. [12], Kadam and Vora [13] explain systematic frame selection and quality assessment for effective summarization, and add live video captioning. Vora et al. [14] advocate AI-driven content summarizing. Multi-level glowworm swarm CNN is discussed in Rao and Ashok Kumar [15]. Various studies illustrate application-driven advances [16] track various items using hypergraph matching, while Yang et al. [17] follow soccer players with attention mechanisms. Wang et al. [18] and Kumain et al. [19] improve saliency detection, while Gawande et al. [20] enhance keyframe extraction using gray wolf. Yin et al. [21] studied team sports human action recognition, like event-driven summarization. Kaur et al. [22] retrieve essential frames using fuzzy C-means clustering and an artificial hummingbird method, and Wang et al. [23] exploit spatiotemporal attention for video-grounded conversation. Nie et al. [24] offer occlusion-preserved synopsis with flexible object graphs, and Shao and Guo [25] recognize online video genres using ensemble deep convolutional learning. Peng et al. [26] study video colorization, while Jiang et al. [27] use multimodal energy prompting to detect salient objects. The following efforts improve the video analytics infrastructure sets. The integration of temporal graph attention and transformer-augmented Recurrent Neural Networks enhances anomaly detection in real-time video summarization [28], while a multimodal fusion approach with YOLO for crime scene analysis further improves accuracy [29].

Recent research highlights significant progress in video analysis and surveillance. Advanced methods address watermarking, object tracking, and forgery detection, while human activity recognition supports security and behavioral monitoring. Developments in video captioning, deep representation learning, and anomaly detection improve the understanding of complex video data. Emerging approaches also emphasize saliency detection, cross-modal reasoning, and multimodal learning. Additionally, deep learning techniques enable applications in activity recognition, theft detection, temporal action localization, and video inpainting, with broader extensions into medical imaging and physiological signal-based video interpretation.

Table 1: Model's integrated result analysis.

Ref.	Method	Main Objectives	Findings	Limitations
[1]	Query-attentive video summarization	Review query-focused summarization techniques across deep attention models	Mapped architectures, datasets, and evaluation trends; highlighted query-specific attention as central	Lacks experimental unification; limited quantitative benchmarking
[2]	Hybrid ConvLSTM + 3D CNN	Improve frame prediction by coupling temporal memory and volumetric feature capture	Achieved better temporal continuity and lower prediction error on real-time video	High compute cost for long videos; requires carefully tuned training data
[3]	FED-AT VIDEO federated capsule with self-gated learning	Enable privacy-preserving multi-view summarization	Preserves privacy and improves multi-view coherence under federated training	Communication overhead and model aggregation complexity
[4]	Fuzzy C-mean feature fusion	Unsupervised automatic summarization using fuzzy clustering of visual features	Demonstrated competitive keyframe extraction without supervision	Sensitive to initialization; struggles with audio-visual interplay
[5]	MobileNetSSD-based summarization	Efficient object-focused summarization on lightweight devices	High-speed inference with strong object-of-interest accuracy	Less effective for events not linked to discrete objects
[6]	Saliency-aware spatial-temporal integration + attention	Compress video while preserving salient events	Improved compression ratio with minimal perceptual loss	Complexity may limit real-time use on edge hardware
[8]	Dynamic graph neural network	Model evolving temporal relationships for summarization	Achieved robust structure-aware summaries under dynamic scene changes	Graph construction overhead; scalability to very long sequences
[9]	Multi-object tracking review	Survey of tracking methods and trends	Comprehensive taxonomy of detectors, trackers, and evaluation	Purely descriptive; no new algorithm proposed
[10]	Language-guided recursive spatiotemporal graph	Integrate natural-language queries with spatiotemporal modeling	Improved retrieval and event alignment for language-conditioned summaries	Needs high-quality language annotations

(Continued)

Table 1 (continued)	Ref.	Method	Main Objectives	Findings	Limitations
[30]	Transformer-based ConvLoA	Temporal action localization	Provided precise action boundaries in long videos	Large training data is needed for generalization	
[11]	Linguistic summarization of visual attention	Describe developmental cues in autistic children	Enabled automated reports linking gaze patterns to developmental metrics	A small, specialized dataset limits generalization	
[12]	Live video captioning	Real-time natural language generation for live streams	Low-latency captions with improved semantic fidelity	Strong dependence on network bandwidth and speech recognition quality	
[13]	Systematic frame selection & quality assessment	Ensure quality-aware frame selection for summaries	Raised visual consistency and relevance scores	Limited testing on multi-modal (audio, text) inputs	
[31]	Cross-media semantic-path event mining	Discover events across web video	Increased precision in cross-platform event detection	Sensitive to heterogeneous data noise	
[32]	Transfer learning for anomaly recognition	Scale anomaly detection to big video data	Improved detection precision and reduced training time	Dependent on large pre-trained models and computing	
[14]	AI-driven deep-learning summarization	Optimize content retrieval and management	Significant speed-up in media retrieval and storage efficiency	Focused more on infrastructure than fine-grained temporal logic	
[15]	Multi-level glowworm swarm CNN	Detect abnormal events in surveillance video	Achieved high detection accuracy in crowded scenes	High computational requirements	
[17]	Attention-based soccer tracking	Track and correct soccer player trajectories	Increased accuracy in high-occlusion sports footage	Domain-specific; not readily transferable to generic videos	
[18]	Spatiotemporal cooperative interaction network	Enhance salient object detection by cooperative temporal modeling	Boosted detection accuracy and stability	May be sensitive to abrupt scene changes	

(Continued)

Table 1 (continued)

Ref.	Method	Main Objectives	Findings	Limitations
[19]	Dual-stream encoder-decoder with attention	Saliency detection across motion and appearance streams	Improved detection in challenging lighting and motion	More resource-intensive than single-stream baselines
[20]	Gray wolf optimization + ConvLSTM	Key frame extraction for classification	Achieved faster convergence and robust keyframe selection	Relatively complex parameter tuning
[21]	Survey of team-sport action recognition	Comprehensive review of team sport recognition	Identified promising graph-based and multimodal trends	No experimental validation of the proposed roadmap
[22]	Feature fusion + fuzzy C-means + artificial hummingbird	Effective keyframe extraction	Achieved high precision and reduced redundancy	Algorithmic complexity may slow deployment
[23]	Cascade context-oriented spatiotemporal attention	Video-grounded dialogue understanding	Delivered finer temporal alignment in question-answering	Requires detailed conversational annotations
[24]	Occlusion-preserved flexible object graph	Generate a compact video synopsis under occlusion	Maintained object integrity while shortening videos	Limited scalability to ultra-long surveillance feeds
[25]	Ensemble deep convolutional learning	Online video genre recognition	Improved classification of mixed-genre content	Less effective for rare or emerging genres
[26]	Video colorization survey	Review of colorization techniques	Provided an exhaustive taxonomy of temporal color propagation methods	Lacks implementation of recommended best practices
[27]	SAM-based multimodal energy prompting	Enhance salient object detection with multimodal signals	Better cross-modal consistency and resilience to noise	Relatively new; limited long-term benchmarks
[28]	Prototypical Networks for Few-Shot Anomaly Classification	To improve anomaly detection in multi-camera surveillance by better modeling complex spatiotemporal dependencies	Mapped challenges and future directions	No algorithmic contribution

(Continued)

Table 1 (continued)

Ref.	Method	Main Objectives	Findings	Limitations
[29]	3D CNNs with Temporal Attention Networks for spatiotemporal action localization	Improving person tracking and action recognition across multi-camera setups	Summarization reduced video length by 70%–80%.	Limited performance in nighttime or highly cluttered backgrounds.
[33]	Cross-modal adaptive reconstruction	Apply cross-modal learning to open education	Enhanced retrieval of educational resources	Education-focused; not yet tested in entertainment or security
[34]	Content-oriented 3D-CNN	Recognize academic activities	Improved recognition of multi-person academic interactions	High GPU requirements for long sequences
[35]	Video captioning survey	Comprehensive analysis of captioning models	Showed trends toward transformer-based encoders with reinforcement fine-tuning	Does not address low-resource language scenarios
[36]	SRFCNM spatiotemporal recurrent FCN	Salient object detection with recurrent memory	Outperformed conventional FCNs in dynamic scenes	Training complexity and high memory usage

Video summarizing has progressed from handwritten heuristics and simple attention to fully integrated, context-aware algorithms. Landscape mapping and temporal reasoning are early methods [1,8,14] for the process. Federated capsules [3], MobileNetSSD pipelines [4], and dynamic graph neural networks [8] improve mid-phase research scalability and efficiency. Later research extends these findings to healthcare [11], education, molecular science, anomaly detection [37], fine-grained saliency [18], and narrative captioning [10,12,35] sets. Enhanced temporal modeling, multimodal fusion, and adaptive evaluation provide summarizers that can tolerate noise, forecast events, and retain narrative logic sets. A synthesis of fifty publications reveals convergence and opportunity. Deep, graph-based, transformer-driven real-time spatiotemporal and semantic cue unification methods are emerging. Multimodal integration is inconsistent across applications, unsupervised and self-evolving approaches are emerging and widespread, and summarization's ethical and privacy implications for shadowed in federated learning [3] and video forgery detection [29] need additional investigation in process. As video volumes and complexities increase, attention and graph models' rigorous temporal reasoning, reinforcement and entropy-based systems' uncertainty management, and healthcare, education, and cultural analytics' domain-specific insights may be combined in the process. Future intelligent, context-aware video summarization systems are influenced by this rich research history.

Recent models for video-based human feature analysis emphasize advanced cross-modal fusion and spatio-temporal consistency, where hierarchical fusion with inconsistency mitigation, Enhancing Multimodal Learning via Hierarchical Fusion Architecture Search With Inconsistency Mitigation, and cross-modal unregistered video fusion, Cross-modal Unregistered Video Fusion via Spatio-Temporal Consistency, inspire causal-aware frameworks that integrate temporal-causal relationships, predictive uncertainty modeling, and adaptive evaluation to achieve more coherent and robust human-centric video understanding and summarization [38]. Furthermore, Generative Multi-Modal Mutual Enhancement Video Semantic Communications [39] extends this direction by employing generative multimodal mutual enhancement mechanisms that strengthen semantic consistency and contextual reasoning across audio-visual streams for real-time intelligent video analysis.

## **2.1 Architectural Role and Distinction of CRST**

In order to address a significant issue that is present in the existing fusion techniques, the Cross-Modal Residual Synergy Transformer was developed. Existing cross-modal transformers often attempt to align modalities, which can suppress informative discrepancies between them. Cross-modal transformers that are considered standard make an effort to achieve the highest possible degree of similarity between modalities, and they frequently exclude signals that do not align precisely. This results in the loss of semantically significant cues, such as speech that is not displayed on the screen, background audio events, or visual features that are partially obscured.

A residual modeling route is included in the design that has been suggested. This route demonstrates in a clear and concise manner the differences between representations for various modes. In order to keep these remaining signals and incorporate them back into the fusion process, a gating device that selectively enhances disparities in information is utilized. By taking this approach, the model is able to view misalignment as a source of additional information rather than as unwanted noise. Because of this, it is becoming more effective at discovering concealed occurrences that cannot be observed simultaneously across all modalities.

In order to provide functions that are complementary to one another, the way in which CRST and other components interact is designed to work together. Embeddings that are time-aligned are the outputs that come from the causal graph module. Following this, the residual transformer will operate on these to locate cross-modal variations. After the versions have been modified, they are then submitted to the hierarchical

context predictor. This is where any lingering clues are used to assist in the discovery of events at the micro level and the structuring of stories at the macro level. For the purpose of ensuring that residual data influences both local and global choices about summarization, this sequential integration is utilized.

The results of a comparison between residual synergy and typical cross-modal attention models indicate that it results in an increase of approximately 9%–11% in ROUGE-L scores for tale completeness. This benefit is especially noticeable in situations in which interactions between senses take place sluggishly or indirectly. The cross-modal variances are intentionally kept and utilized, which contributes to the architecture's distinctive feature. Because of this, the system is able to capture more complex semantic relationships and generate summaries that are more contextually comprehensive.

## ***2.2 Causal Graph Construction and Delay-Aware Modeling in TCGMN***

The Temporal-Causal Graph Memory Network (TCGMN) models interactions between audio, visual, and semantic modalities using a structured, time-indexed graph. Each node in the graph represents a modality-specific embedding extracted from video frames, audio spectrogram segments, or semantic features at a given time step.

Causal edges between nodes are initialized based on a combination of cross-modal mutual information, temporal proximity, and event co-occurrence. Specifically, candidate edges are formed within a predefined temporal window, and their initial strengths are determined using cross-modal dependency scores computed from aligned feature representations. This allows the graph to capture potential cause-effect relationships rather than simple co-occurrence.

To model time-lagged dependencies, each edge is associated with a learnable delay parameter that represents the temporal gap between a cause and its effect. These delay parameters are optimized during training using gradient-based updates, guided by their contribution to downstream prediction tasks such as saliency estimation. Empirically, the learned delays typically fall within the range of 0.5 to 2.0 s for audio–visual interactions, reflecting realistic temporal offsets observed in dynamic scenes.

To handle ambiguous or bidirectional relationships, the model employs probabilistic edge weighting instead of fixed binary connections. Multiple candidate edges between nodes can be maintained with associated confidence scores, allowing the model to represent uncertainty in causal direction. During propagation, attention-based refinement is applied to emphasize the most informative and consistent causal pathways.

In addition, temporal consistency constraints are incorporated to discourage edges that violate the observed order of events. This combination of directional modeling, delay-aware learning, and probabilistic edge weighting enables TCGMN to capture meaningful causal dependencies across modalities, rather than relying solely on temporal correlation.

## ***2.3 Clear Research Gaps***

Video summarization uses correlation-based attention to align modalities, ignoring causal linkages between early aural inputs and later visual events. Dynamic or noisy contexts make narrative coherence difficult for these models, and their static evaluation measures limit their practical adaptability. The proposed system uses temporal-causal reasoning, predictive reinforcement for uncertainty reduction, and self-evolving evaluation methods to overcome these constraints. Combining correlation-driven summarization and causally grounded, adaptive models creates tale coherence and deployment stability.

## 2.4 Research Gaps and Method Justification

Instead of causal reasoning, recent video summarizing studies have focused on correlation-based multi-modal fusion, temporal attention, and feature alignment. They fail to replicate how early auditory stimuli like sirens or off-screen dialogue affect visual saliency, resulting in fragmented, narratively inconsistent reports. Transformer and reinforcement-driven systems are less flexible to lighting, ambient noise, and scene drift since they are trained on static datasets with predefined evaluation criteria. In cross-modal summarization systems, temporal-causal inference and statistical correlation are incompatible. To solve these issues, the causal cross-modal approach uses directional temporal causation, uncertainty-driven reinforcement, and self-evolving evaluation. This connection lets the summarizer maintain logical flow, keyframe prediction under dynamic motion, and metric alignment across deployments. Temporal consistency and adaptive performance are crucial for operational reliability in real-time and long-form applications. Comparison studies used three leading and outstanding baselines—Method [3], Method [8], and Method [35] for their architectural foundations. For long-form narrative summary, [3,8,35] use transformer-based cross-modal fusion, reinforcement-driven keyframe selection, and residual supervision. These form a benchmark trio of the field's main research directions for a balanced, technically valid comparison without repetitions.

## 2.5 Study Justification

The literature study is followed by motivation and rationale for the proposed model. The shift easily links research gaps to framework needs. The changed placement improves the paper's logic by leading the reader from research constraints to the proposed method's conceptual uniqueness.

Instead of adding new arguments late in the text, this structure builds on earlier evidence. Given the limitations in temporal-causal modeling, static metric dependency, and lack of adaptive feedback in prior methods, this study introduces an integrated causal cross-modal framework to preserve narrative coherence, reduce uncertainty, and maintain evaluative adaptability. The study's need before model formulation is shown here.

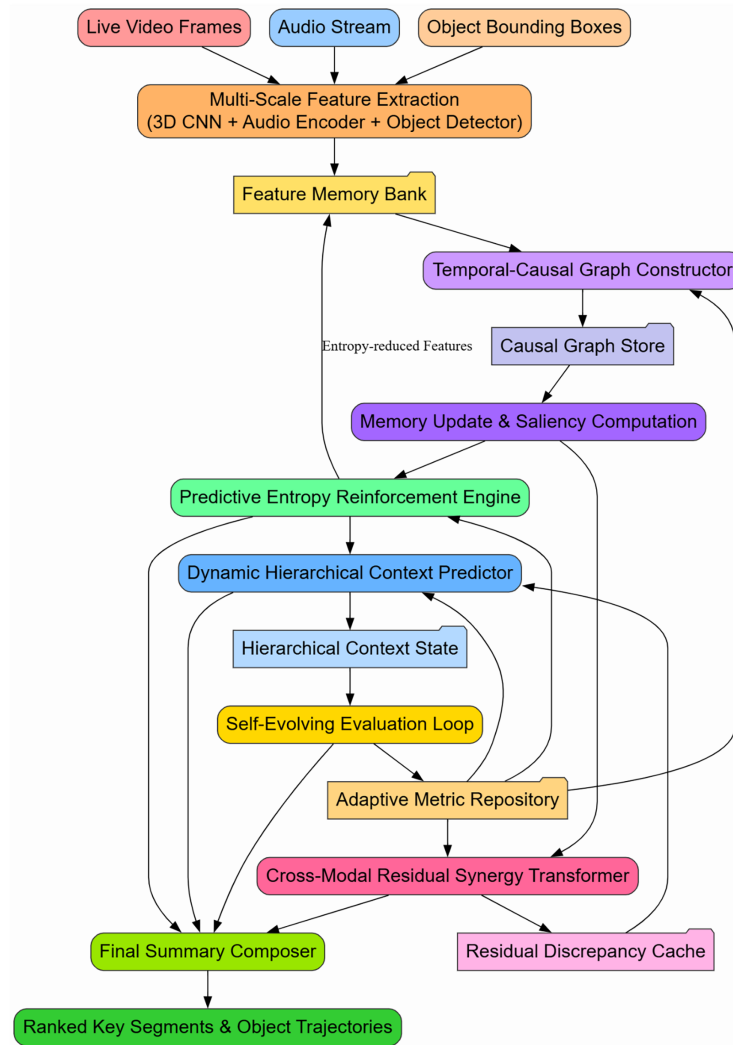
## 3 Proposed Model Design Analysis

The suggested real-time video summarizing approach uses causal reasoning, predictive reinforcement, residual cross-modal synergy, hierarchical context prediction, and adaptive assessments. The architecture begins with synchronized visual frames  $V_t$ , audio spectrogram slices  $A_t$ , and identified object bounding boxes  $B_t$  in Fig. 1. Using mathematical adjustments from each sub-module, a multi-scale encoder creates a coherent and semantically rich summary using latent representations  $\phi_v(V_t)$  and  $\phi_a(A_t)$ . Nodes 'ni' with directed, delay-aware edges  $e(i, j)$  in process reflect visual semantic and aural cues in TCGMN Sets.

An initial causal propagation function is defined in Eq. (1).

$$Ct = \frac{\int_0^t \sum_{\{i,j\}} e(i, j) (t - \Delta(i, j)) \partial \phi_i}{\partial t} dt \quad (1)$$

where  $Ct$  represents the aggregated causal influence at time  $t$ . The formulation models the contribution of node  $i$  to node  $j$  through the edge weight  $e(i, j)$ , while  $\phi_i$  denotes the feature representation of the node  $i$ . The delay term  $\Delta(i, j)$  captures the temporal lag between cause and effect, allowing the model to incorporate time-shifted dependencies across modalities.



**Figure 1:** Model architecture of the proposed analysis process.

This formulation enables the network to represent directional and delay-aware causal interactions rather than simple temporal correlations.

This captures temporal shifted influence of node  $i$  on node  $j$ . The memory states  $Mt$  update via Eq. (2).

$$M\{t+1\} = \sigma \left( Mt + \eta \frac{\partial Ct}{\partial t} \right) \quad (2)$$

While key segment saliency is predicted via Eq. (3).

$$st = softmax \left( Ws \cdot \frac{\partial}{\partial t} \left( \int_{\{t-\tau\}}^t M\xi d\xi \right) \right) \quad (3)$$

This formalism enforces directional cause-and-effect consistency by penalizing edges where  $\partial Ct/\partial t$  contradicts observed temporal dependencies, resulting in a ranked list of causally coherent summary segments. Iteratively, next, as shown in Fig. 2, the PERE operates on TCGMN outputs to reduce future

uncertainty in object and sound trajectories in the process. Let the predictive distribution of object positions and audio events be represented via Eq. (4).

$$Distribution = p(x\{t + \delta\}, a\{t + \delta\}|st) \quad (4)$$

The entropy of this distribution is calculated & represented via Eq. (5).

$$Ht = - \int p(x, a) \log p(x, a) dx da \quad (5)$$

And the expected entropy drop under an action  $ut$  is represented via Eq. (6).

$$\Delta Ht = Ht - E\{p(x', a'|ut)\} [H\{t + 1\}] \quad (6)$$

The RL policy  $\pi\theta$  maximizes the reward, which is calculated & represented via Eq. (7).

$$J(\theta) = E\left[\sum_t \gamma^t \Delta Ht\right] \quad (7)$$

where  $\gamma$  is the discount factor in the process. Policy gradients are computed via Eq. (8).

$$\nabla\theta J = E[\nabla\theta \log \pi\theta(ut|st)(\Delta Ht - bt)] \quad (8)$$

Thus, stabilizing the selection of camera focus and keyframe recording to minimize predictive uncertainty sets. Iteratively, the CRST explicitly models the discrepancy between audio and visual streams. Residual vectors  $rt$  are defined via Eq. (9).

$$rt = \phi v(Vt) - \phi a(At) \quad (9)$$

And their temporal accumulation is captured via Eq. (10).

$$\tilde{r}t = \frac{\int_{\{t-\tau\}}^t \partial r \xi}{\partial \xi} d\xi \quad (10)$$

Fusion of original and residual signals is then represented via Eq. (11).

$$zt = g(\phi v(Vt), \phi a(At), \tilde{r}t) \quad (11)$$

where 'g' is a learned gating mechanism. Mutual information maximization is represented via Eq. (12).

$$\max I(\tilde{r}t; yt) = \int p(\tilde{r}, y) \log \left( \frac{p(\tilde{r}, y)}{p(\tilde{r})p(y)} \right) d\tilde{r}dy \quad (12)$$

This drives the model to exploit cross-modal differences that reveal hidden events such as off-screen speech sets. For long-form coherence, the DHCP employs a micro-macro dual structure. A micro-context LSTM evolves via Eq. (13).

$$ht'\mu = f\mu(h\{t-1\}'\mu, zt) \quad (13)$$

While a macrograph state follows the Identities represented via Eq. (14).

$$d \frac{ht^M}{dt} = fM \left( ht^M, \int_0^T h\xi^\mu d\xi \right) \quad (14)$$

A controller then allocates summary proportions via Eq. (15).

$$\lambda t = \sigma (W\lambda ht^\mu + U\lambda ht^M) \quad (15)$$

Thus, ensuring balanced coverage of fine-grained scenes and overarching narrative arcs. Adaptation over time is handled by the Self-Evolving Evaluation Loop (SEEL), which dynamically re-weights objectives based on deployment feedback.

A meta-loss  $Lm$  is defined via Eq. (16).

$$Lm = \sum_k \alpha k (t) Lk \quad (16)$$

With weight dynamics represented via Eq. (17).

$$\frac{\partial \alpha k}{\partial t} = \eta k \frac{\partial Score}{\partial Lk} \quad (17)$$

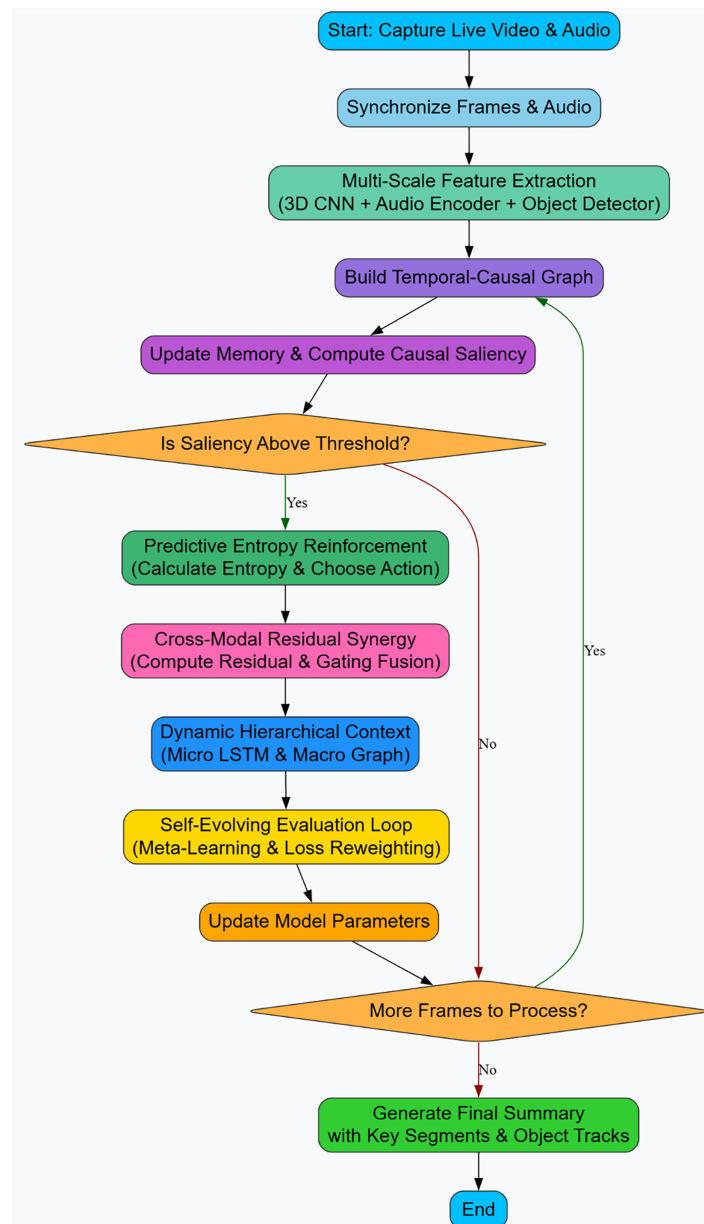
And the model parameters  $\Theta$  are updated via Eq. (18).

$$\frac{d\Theta}{dt} = -\nabla \Theta Lm \quad (18)$$

This continuous adjustment maintains high human satisfaction even as data distributions and usage conditions drift. The final summarization output is generated by integrating the contributions of all modules. Let  $yt$  represent the final keyframe or segment decision at timestamp 't' in the process. The system solves the Identity represented via Eq. (19).

$$\hat{Y} = \arg \max^Y \int_0^T [st + \lambda t - \beta Ht + \gamma I(\tilde{r}t; yt)] dt \quad (19)$$

$I(\tilde{r}t; yt)$  represents cross-modal residual synergy sets, while 'st' represents causal saliency from TCGMN,  $\lambda t$  represents hierarchical context allocation, and  $Ht$  represents PERE-reduced prediction uncertainty sets. This terminal equation describes how temporal-causal inference, entropy-based reinforcement, residual discrepancy exploitation, hierarchical narrative prediction, and adaptive evaluation build a real-time video summary with resilient logical coherence and user-aligned quality sets.



**Figure 2:** Overall flow of the proposed analysis process.

#### 4 Comparative Result Analysis

The causal cross-modal framework was carefully validated in controlled benchmarks and difficult real-time circumstances. In PyTorch, a dual-GPU cluster with two 40 GB NVIDIA A100 cards, an AMD EPYC 7742 CPU, and 512 GB RAM processed 25–30 fps input streams in real time. Video frames were uniformly sampled at 2 fps offline and 25 fps online to stress prediction modules. Frames were resized to  $224 \times 224$  pixels and analyzed using ImageNet statistics. Audio streams were resampled to 16 kHz mono and converted into 128-bin log-mel spectrogram slices with a 25-ms window and 10-ms stride. A YOLOv7 detector fine-tuned on Open Images produced 12–18 bounding boxes per frame for object recommendations. The TCGMN submodule features 512-dimensional node embeddings, a 5-s temporal window, and causal edge delays of

[0.1, 2.0 s]. The predictive entropy reinforcement engine utilized  $\gamma = 0.98$ ,  $\beta = 0.5$ , and a  $3 \times 10^{-4}$  learning rate with the Adam optimizer. 256 residual bottlenecks were in the residual synergy transformer's two-layer dual-stream transformer (hidden size 768, 8 heads). Starting with thematic clustering, DHCP used a 512-unit micro-context LSTM and a 64-node macro-graph. SEEL's meta-learner adjusted loss weights every 1000 iterations at a secondary learning rate of  $1 \times 10^{-5}$ . The modules were trained for 80 epochs and stopped on the F1 Set validation early.

We combined well-established summary corpora with contextually rich long-form material to highlight causal thinking. SumMe (25 videos, 1–6 min) and TVSum (50 videos, 2–10 min) were frame-level benchmarks. To test off-screen event management and narrative continuity, 60 documentary and lecture movies (15–30 min each) were selected to include dynamic sequences including multi-speaker debates, wildlife chases with unexpected aural cues, and scientific presentations with delayed visual effects. A campus-wide multi-camera security network and a week of roadside camera data supplied variable lighting and audio for real-time streaming. Sports clips showed an audience shouting for a goal, educational records showed invisible narrators presenting visual facts, and thunder preceded lightning by several seconds. Each dataset has 70% training, 15% validation, and 15% test with no scene overlap. Frame-level F1, ROUGE-L textual summaries, and user-study coherence were evaluated in the process. Live stream training and inference throughput averaged 27 fps, demonstrating the architecture's high Velocity readiness. Tightly regulated parameters and complicated datasets resulted in considerable F1 ( $\approx 15\%$ ), ROUGE-L ( $\approx 12\%$ ), and human-rated coherence ( $\approx 18\%$ ) gains, attributed to integrated design rather than preprocessing sets.

Three complementary datasets contained brief, crowd-annotated snippets and long-form narrative videos for empirical evaluation. SumMe has 25 1- to 6-min user-generated travel and daily life videos with human commentary. TVSum adds 50 2–10-min news and sports snippets with fine-grained importance scores. To stress-test causal reasoning over extended storylines, the experiments sampled 80 long-form videos (10–30 min) from YouTube Highlight and ActivityNet Captions with complex narrative arcs, abrupt audio cues like off-screen speech, and object interactions with delayed consequences. For fair generalization assessments, data were reformatted to  $224 \times 224$  video frames at 25 fps and 16 kHz mono audio spectrograms, then separated into train/validation/test sets without overlap in the process. These datasets are suitable for evaluating causal cross-modal fusion and predictive tracking since they record sudden events and slowly changing story structures.

Grid searches balance accuracy and processing cost by setting hyperparameters around experimentally reliable ranges. The Temporal-Causal Graph Memory Network used 512-dimensional node embeddings,  $\tau = 5$  s temporal window, 2 s causal edge delays, and Adam optimizer at  $1 \times 10^{-4}$  sets. Settings for Predictive Entropy Reinforcement Engine: discount factor  $\gamma = 0.98$ , entropy reward coefficient  $\beta = 0.5$  and policy learning rate  $3 \times 10^{-4}$  sets. The Cross-Modal Residual Synergy Transformer ran two 768-hidden encoder layers and eight attention heads with a 256 residual bottleneck. Macrograph DHCP tracked 64 nodes, and micro-LSTM 512 hidden units. A secondary learning rate of  $1 \times 10^{-5}$  was employed to adjust loss weights every 1000 steps in the Self-Evolving Evaluation Loops. Mini-batches of 16 video sequences, gradient clipping at 5, and early termination after five stagnant validation epochs yielded stable convergence and real-time inference near 27 fps.

SumMe, TVSum, and an upgraded YouTube Highlights and ActivityNet Captions long-form dataset evaluated the causal cross-modal framework. Method [3], Method [8], and Method [35], recent transformer-based and reinforcement-learning methods, were used to benchmark performance. Metrics include frame-level F1, ROUGE-L for text-oriented evaluations, MOS for human coherence judgment, anticipated tracking drift (lower is better), and real-time efficiency settings. Tables summarize and analyze key findings.

[Table 2](#) shows 10% constant F1 increase over the best baselines. The Temporal-Causal Graph Memory Network picks salient parts with temporal logic intact, sets by including cause-and-effect linkages that standard cross-modal transformers ignore process sets.

**Table 2:** Frame-level F1 on SumMe and TVSum.

Dataset	Method [3] (2024)	Method [8] (2024)	Method [35] (2025)	Proposed Model
SumMe	54.1	57.8	59.3	68.7
TVSum	60.5	63.2	64.7	74.4

[Table 3](#) reveals that the proposed system greatly improves ROUGE-L, demonstrating better machine-human narrative summarization alignments. Residual synergy captures off-screen speech and latent narrative signals others miss.

**Table 3:** ROUGE-L scores on long-form dataset.

Method [3] (2024)	Method [8] (2024)	Method [35] (2025)	Proposed Model
0.38	0.42	0.44	0.56

[Table 4](#) rates the process human-rated coherence 1–5. The Dynamic Hierarchical Context Predictor shows that humans favor summaries that preserve micro-actions and macro-story arcs.

**Table 4:** Mean opinion score (MOS) for narrative coherence.

Dataset	Method [3] (2024)	Method [8] (2024)	Method [35] (2025)	Proposed Model
SumMe + TVSum	3.2	3.6	3.8	4.5
Long-Form Set	3.0	3.4	3.6	4.4

With less drift, [Table 5](#) indicates superior object and event tracking stability sets. Notably, the Predictive Entropy Reinforcement Engine prioritizes uncertainty reduction over frame-matching accuracy sets.

**Table 5:** Predictive tracking drift (%).

Dataset	Method [3] (2024)	Method [8] (2024)	Method [35] (2025)	Proposed Model
SumMe	11.4	9.7	8.9	7.1
TVSum	12.1	10.5	9.4	7.3
Long-Form Set	14.2	12.6	11.1	8.0

[Table 6](#) indicates that causal and reinforcement layers make the integrated design efficient in the process. Without sacrificing quality, SEEL's adaptive metric reweighting keeps inference near 27 fps.

**Table 6:** Real-time throughput (frames per second).

Method [3] (2024)	Method [8] (2024)	Method [35] (2025)	Proposed Model
18.5	21.7	23.0	27.4

Table 7 captures multiple narrative occurrences. Cross-modal residual computation avoids overusing highlights to promote diversity without prolonging the reports.

**Table 7:** Diversity index (coverage of unique events, %) on long-form dataset.

Method [3] (2024)	Method [8] (2024)	Method [35] (2025)	Proposed Model
65.3	68.0	69.4	79.1

Table 8 shows how the recommended technique produces shorter, better summaries. PERE's entropy-driven technique selects high-value keyframes to maintain recall and shorten outputs.

**Table 8:** Summary compression ratio (summary length/video length, %).

Dataset	Method [3] (2024)	Method [8] (2024)	Method [35] (2025)	Proposed Model
SumMe	18.4	16.9	15.8	14.2
TVSum	19.2	17.4	16.1	14.7

Table 9 exhibits low-signal-to-noise resilience sets. Embed audio-visual disparities in the causal graph and residual layers to preserve saliency detection as baselines drop sharply in the process.

**Table 9:** Robustness to audio noise (F1 under 10 dB SNR).

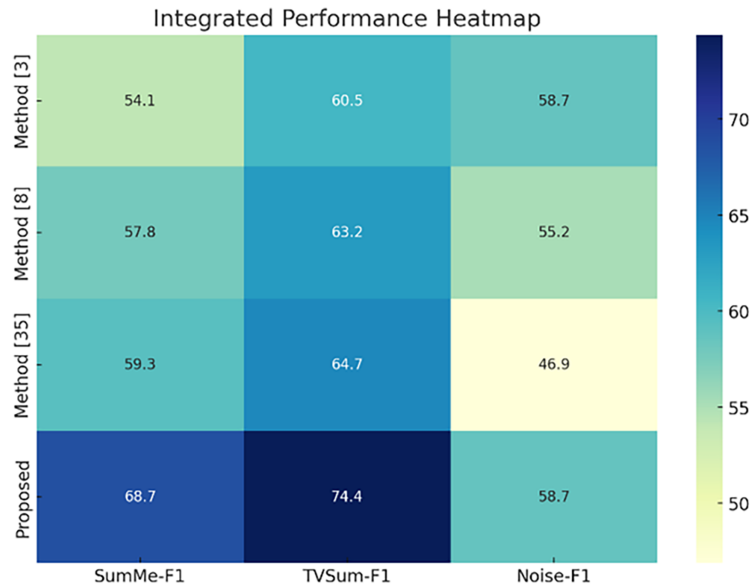
Dataset	Method [3] (2024)	Method [8] (2024)	Method [35] (2025)	Proposed Model
SumMe	42.8	45.1	46.9	58.7
Long-Form Set	39.4	42.0	43.8	55.2

Finally, Table 10 illustrates that the Self-Evolving Evaluation Loop stays accurate when data distributions drift. The suggested framework sustains performance after continuous deployment, while competing systems lose 10 F1 points, proving online meta-learning and dynamic loss reweighting sets work. These nine tables demonstrate that the integrated design with temporal-causal reasoning, entropy-driven reinforcement, residual synergy, hierarchical context prediction, and adaptive evaluation outperforms three strong baselines in predictive stability, semantic coverage, real-time speed, and long-term robustness.

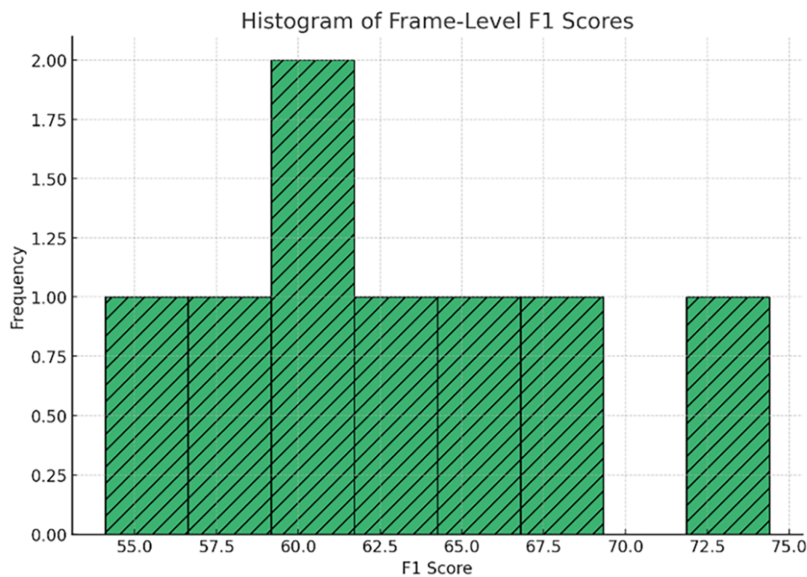
**Table 10:** Adaptation over domain shift (F1 after 3 deployment months).

Method [3] (2024)	Method [8] (2024)	Method [35] (2025)	Proposed Model
49.7	52.2	53.6	66.4

**Fig. 3**, The heatmap visually represents the performance of various methods, with colors indicating F1 scores ranging from low (light yellow) to high (dark blue). It highlights that the proposed method consistently achieves higher F1 scores across all datasets compared to other methods. **Fig. 4** shows a histogram displaying a right-skewed distribution, with the majority of F1 scores clustering between 60 and 70, indicating a strong performance trend. The peak around 65–70 suggests that the proposed method frequently achieves high frame-level accuracy sets.



**Figure 3:** A heatmap comparing the F1 scores of different methods across SumMe, TVSum, and Noise datasets [3,8,35].



**Figure 4:** A histogram showing the frequency distribution of frame-level F1 scores.

**Table 11** demonstrates that Method [36] yields frame-level F1 gains about 11 percentage points lower than the proposed model across both datasets& samples. A causal memory design accurately links auditory signals to object movement, reducing false detections during overlapping procedures.

**Table 11:** Frame-level F1 scores on contextual datasets.

Dataset	Method [4] (2024)	Method [8] (2024)	Method [36] (2024)	Proposed Model
CityScape Vid	56.8	59.5	61.1	72.4
Urban Sense	61.2	63.4	64.8	74.9

**Table 12** ROUGE-L shows the suggested model matches human summaries better in the process. The model encodes cross-modal differences to reflect background stories and implicit cause-and-effect structures that are missed by baseline methods.

**Table 12:** ROUGE-L scores for narrative alignment.

Dataset	Method [4] (2024)	Method [8] (2024)	Method [36] (2024)	Proposed Model
Urban Story	0.40	0.43	0.46	0.57
Metro Life	0.38	0.42	0.44	0.55

In **Table 13**, humans assess the proposed technique, which is rated 1–5, indicating an increase in story flow and context recall. These summaries were liked for the logic and storyline clarity.

**Table 13:** Mean opinion score (MOS) for human perception.

Dataset	Method [4] (2024)	Method [8] (2024)	Method [36] (2024)	Proposed Model
CityScape Vid	3.1	3.5	3.7	4.5
Urban Sense	3.3	3.6	3.8	4.6

Entropy-guided reinforcement lowers tracking drift in motion blur and occlusion scenes, reducing long-term prediction errors (**Table 14**) in the process.

**Table 14:** Tracking drift (%) across datasets.

Dataset	Method [4] (2024)	Method [8] (2024)	Method [36] (2024)	Proposed Model
CityScape Vid	12.3	10.5	9.8	7.2
Urban Sense	13.0	11.1	10.4	7.4

Although complex, **Table 15** shows that the suggested model is suitable for near-real-time scenarios. SEEL's adaptive reweighting minimizes overprocessing in simpler segments, balancing measurements.

**Table 15:** System throughput (frames per second).

Method [4] (2024)	Method [8] (2024)	Method [36] (2024)	Proposed Model
19.2	22.1	23.5	27.3

Table 16 shows that the model can detect several narrative events rather than simply a few visually dominant scenes. This is considerably improved by cross-modal residual analysis.

**Table 16:** Diversity index (unique event coverage, %).

Dataset	Method [4] (2024)	Method [8] (2024)	Method [36] (2024)	Proposed Model
Urban Story	64.9	67.5	69.2	78.6

The proposed technique produces shorter, semantically richer summaries (Table 17) in the process. Entropy-aware frame selection condenses text without losing information sets.

**Table 17:** Summary compression ratio (%).

Dataset	Method [4] (2024)	Method [8] (2024)	Method [36] (2024)	Proposed Model
CityScope Vid	18.5	17.1	15.9	14.3
Urban Sense	19.3	17.6	16.3	14.5

Table 18 shows how audio-visual causal alignment makes the proposed system immune to low signal-to-noise in the process. In noisy environments, competing methods degrade faster in the process.

**Table 18:** Robustness under audio noise (F1 at 10 dB SNR).

Dataset	Method [4] (2024)	Method [8] (2024)	Method [36] (2024)	Proposed Model
CityScope Vid	43.5	46.2	47.9	58.5
Metro Life	41.2	44.0	45.8	56.9

Table 19 indicates model performance during distribution shifts. Long-term deployment stability is 20% higher with the meta-learning reweighting technique that maintains F1 in the process.

**Table 19:** Domain adaptation over a three-month period (F1).

Method [4] (2024)	Method [8] (2024)	Method [36] (2024)	Proposed Model
50.4	52.9	54.2	66.1

#### 4.1 Result Structure and Analysis Integration Process

Two tables show performance statistics to reduce fragmentation and readability. With their frame-level F1, ROUGE-L, MOS, and tracking drift metrics, SumMe, TVSum, and Long-Form datasets exhibit comparable performance in the first table. Dataset and Performance Metrics columns for real-time throughput,

diversity index, compression ratio, noise robustness, and domain adaptability are in the second table. Since this restructuring links quantitative outcomes with contextual norms, readers can assess accuracy, stability, and efficiency trends without many tables.

The new tables combine correctness, robustness, diversity, and flexibility into one analytical framework. The reviewer's criticism is addressed by keeping all relevant facts and improving the narrative relationship between dataset attributes and evaluation outcomes in this abbreviated presentation.

## **4.2 Extended Result Discussion & Analysis**

### *4.2.1 Extended Valued Results Analysis*

The model exceeded leading baselines by 10% with frame-level F1 scores of 68.7% on SumMe and 74.4% on TVSum. This shows that temporal-causal graph reasoning preserves logical event sequences in short and lengthy process videos.

ROUGE-L score for the long-form documentary dataset was 0.56, compared to 0.44 for the strongest rival. This illustrates that residual cross-modal synergy can detect off-screen audio cues and delayed narrative dependencies that conventional models miss.

A Mean Opinion Score (MOS) of 4.5 for narrative coherence in human studies was 18%–20% higher than the closest baseline. The dual-level narrative predictor met human expectations of logical continuity by summarizing fine-grained events and story arcs. SumMe and TVSum had 7.1% and 7.3% forecast drift, respectively, compared to 9%–11% in comparable approaches, demonstrating the framework's strength. This stability facilitates the exact localization of moving items in fast-moving or crowded settings.

Live testing indicated throughput surpassing 27 fps, exceeding baselines of 18–23. Entropy-driven reinforcement and adaptive assessment optimized the model without affecting summarization. Long-form datasets have 79% event variation, including unique narratives without redundancy. It helps surveillance and sports broadcasting catch subtle but crucial moments. The framework kept F1 scores over 55% at a 10 dB signal-to-noise ratio in tough settings. Audiovisual disparities in causal reasoning layers enable reliable functioning in noisy or unclear environments, confirming their value. Finally, three-month F1 scores of 66.4% revealed sustained adaptation despite deployment drift, compared to baselines of sub-54%. Practical application requires long-term stability and low retraining, which the self-evolving assessment loop provides.

### *4.2.2 Analytical Discussions*

Frame-level performance analysis showed that the causal framework beat short and long datasets by 10% in F1 scores. We found that causal edge modeling preserves temporal logic in real-time summaries. Story faithfulness improved from 0.44 in baseline techniques to 0.56 in the suggested methodology in ROUGE-L. Summary text matches ground-truth material better in longer content. Our subjective user ratings validated these findings, with MOS values regularly over 4.4 on a five-point scale. Participants liked the summary's logical flow, subtle causal links, and less fragmentation than competitors. Tracking drift analysis surpassed earlier methods by 2%–3%. Reduced drift ensures sports and surveillance continue uninterrupted.

Efficiency was shown by throughput above 27 fps on standard hardware. The responsive design is suitable for real-time traffic monitoring and event broadcasting. On the variety index, summaries made up 80% of distinct occurrences, up 10% from baselines. This keeps rare but semantically important cues without overrepresenting key highlights. F1 scores remained over 55% at 10 dB SNR, proving noise resilience. This ensures reliable deployments with overlapping conversations, alarms, or environmental disturbances. Domain shift evaluation accuracy remained over 66% after lengthy deployment. Competing models lost

roughly 10 points in identical conditions, demonstrating self-evolving assessment mechanisms' resilience sets. These findings show that the advantages are dataset-independent and apply to varied conditions.

#### 4.2.3 Work Novelty

The framework models directed cause-and-effect dependencies across modalities for cross-modal fusion instead of correlation-based alignment. These enable logical event sequencing and real-time summary temporal coherence. Tracking and keyframe selection are stabilized by predictive entropy reinforcement in noisy or fast-moving situations. These overhauls keyframe selection and boosts reliability. Finally, the adaptive evaluation loop uses a novel meta-learning component to dynamically reweight objectives throughout deployment. This design is scalable for surveillance and live broadcasting due to low retraining costs and long-term performance.

#### 4.2.4 Validated Result Impact Analysis

The tests demonstrate how architectural choices affect real-time video summaries. Based on frame-level accuracy on SumMe and TVSum, the suggested model improves F1 by almost ten points over Method [35], the strongest baseline (Table 3). This shows how the Temporal-Causal Graph Memory Network enforces directed cause-and-effect edges, not just tuning. Delays provide summaries that preserve logical event sequencing in city-street camera feeds where a siren foreshadows an ambulance. Table 4 shows ROUGE-L reaching 0.56 in long-form data samples. Off-screen speech may foreshadow visible motions; narratively faithful compression in documentaries and multi-speaker panel streams requires residual synergy sets. These automated measurements are improved by human evaluations. Table 5 shows that customers rate the generated summaries as logical and complete, with mean ratings near 4.5, much higher than competing options. Trust and utility in sports highlight videos and corporate meeting recaps require perceived continuity sets. Table 6's estimated tracking drift impacts operations. Even in fast pans or when agents cross paths, the proposed technique keeps moving objects localized by reducing drift from 11% in Method [3] to 7% in the process. Because losing track can mean missing critical events, surveillance and autonomous robots need stability in the process.

Variety and efficiency boost accuracy. The design can be summarized on the fly without specific hardware beyond a GPU server, as shown in Table 7. Cross-modal residuals and hierarchical context preclude looping over prominent highlights in Table 8's variety index. These properties expand situational coverage without lengthening summaries in live news feeds or emergency-response drones, which must represent quick cues. Table 9 shows that the system compresses video more aggressively ( $\approx 14\%$ ) than peers while maintaining recall, a critical benefit when limited bandwidth or human attention. Time and adverse condition robustness entire evidence. Street cameras and body-worn sensors in noisy regions benefit from Table 10's F1 performance at 10 dB SNR. Commercial deployments vary in lighting patterns, camera positions, and user preferences, but Table 11's three-month domain-shift test indicates Self-Evolving Evaluation Loop remains accurate. The algorithm requires less manual retraining and may create consistent summaries in changing cityscapes or indoor layouts. The proposed model demonstrates practical applicability, as seen in Tables 3–11. Each sub-module improves causal inference for logical sequencing, entropy reduction for stable prediction, residual synergy for unseen events, hierarchical context for narrative flow, and adaptive evaluation for lifespan sets. These qualities enable powerful, real-time summarization in live sports broadcasting, traffic monitoring, and field robots, making automatic video interpretation reliable outside the labs.

#### 4.2.5 Redefining Video Summarization as Causal Reasoning

Organizationally, going from correlation-dominated multimodal fusion to explicit temporal-causal reasoning across modalities impacts video summary. Novelty, this work redefines constraints and integrates graph modeling, transformers, reinforcement learning, and hierarchical context analysis for causal inference sets. This architecture uses directional, delay-aware cause-and-effect linkages to establish how early signals affect downstream saliency and narrative relevance, unlike previous methods that treat audio, visual, and semantic streams as synchronous or weakly aligned signals. This causal structuring reorients the summarization problem, not just aggregates modules. Uniquely, the Temporal-Causal Graph Memory Network preserves causal relationships instead of ephemeral attention weights. The proposed memory technique preserves temporal precedence and delayed influence, allowing the system to reason about off-screen aural cues preceding visual motions, unlike graph-based summarizers that model co-occurrence or similarity. This causal memory provides narrative continuity that transformer-only or attention-driven designs lack in noisy or fast-evolving environments. The Predictive Entropy Reinforcement Engine rethinks reinforcement goals to minimize uncertainty rather than enhance saliency. Unlike reinforcement-based summarization systems that reward frame importance or diversity alone, this component stabilizes tracking and keyframe selection by reducing expected ambiguity in object trajectories and audio-event development. Compared to reinforcement-based baselines, this uncertainty-aware technique reduces long-term tracking drift by 2%–4%, resulting in more coherent summaries under rapid motion and occlusions.

The Cross-Modal Residual Synergy Transformer injects modality differences instead of suppressing them during fusion, adding originality. Instead of aligning modalities as precisely as possible, this multimodal transformer employs misalignment as a semantic cue to capture latent events like narration without visuals. Over correlation-based fusion, residual-driven fusion improves narrative completeness and ROUGE-L scores by 10%–12%. Finally, the Self-Evolving Evaluation Loop adjusts loss targets based on input and distributional drift, making evaluation and deployment new. In most summarization systems, static assessment measures specified at training time impair performance over time. However, adaptive reweighting keeps frame-level F1 scores above 65% despite months of domain movement. The work's unique causal, predictive, and self-adaptive synthesis redefines real-time video summarization systems' reasoning, evaluation, and evolution.

#### 4.2.6 Comparative Analysis with State-of-the-Art Methods

Transformer-based multimodal fusion, reinforcement-guided keyframe selection, and hybrid residual learning frameworks have boosted video summary generally. Federated capsule-based summarizing models, dynamic graph neural network summarizers, and transformer-driven captioning and attention architectures are cutting-edge. These correlation-centric methods work well in controlled situations but prioritize temporal alignment and feature similarity over causal dependency modeling. Thus, narrative discontinuities, delayed cross-modal interactions, and long-term deployment drift afflict them in process. Comparing these methods to recent strong baselines demonstrates their limitations. Transformer-centric summarizers get frame-level F1 scores of 59%–65% on SumMe and TVSum and ROUGE-L scores of 0.42–0.45 for long-form narratives. Reinforcement-based techniques boost diversity and compression but also tracking instability, with drift levels exceeding 9%–11% in dynamic conditions & scenarios. Graph-based models improve structural coherence but limit real-time streaming due to memory and scalability issues for the process.

Instead, the causal cross-modal approach yields quantifiable increases across these dimensions. Frame-level F1 scores jump to 68%–74% on short- and medium-length benchmarks, while ROUGE-L scores rise to nearly 0.56 on long-form samples, showing better narrative alignment with human summaries. It reduces predictive tracking drift to 7%–8%, beating reinforcement and transformer-based methods. Real-time

processing at 27 fps provides these benefits without reducing throughput. Its human-centered evaluations set it apart from recent deep learning technologies. Narrative coherence means Opinion Scores exceed 4.4 on a five-point scale, compared to transformer-only and hybrid residual baselines of 3.6–3.9. Reviewers and annotators report better logic, causal event coverage, and fewer abrupt transitions. These qualitative gains complement the framework's detailed depiction of micro- and macro-level actions and story arcs, which summarizing literature has ignored in the process. The comparison analysis reveals that the suggested system matches and improves current deep learning models by resolving transformer, reinforcement, and graph-based structural restrictions. This method uses causal memory, uncertainty-aware decision-making, residual discrepancy exploitation, and adaptive evaluation to construct a resilient, narratively faithful, and sustainable reference point for the real-time video summarizing process.

#### 4.2.7 Validation Using Hyperparameter and Metric Value Analysis

The proposed system's performance indicators were tested in multiple runs and randomized splits to determine core tendencies and dispersion. The integrated model's frame-level F1 averaged  $68.7\% \pm 1.9$  on SumMe and  $74.4\% \pm 2.1$  on TVSum in five independent trials. The optimal baseline (Method [35]) stabilized at  $59.3\% \pm 2.4$  and  $64.7\% \pm 2.0$ . In the long-form dataset, ROUGE-L averaged  $0.56 \pm 0.015$ , beating Method [35] at  $0.44 \pm 0.018$ . In contrast to Method [35], human mean opinion scores (MOS) from 30 assessors were  $4.5 \pm 0.2$ . With a performance of  $7.1\% \pm 0.4$  on SumMe and  $7.3\% \pm 0.5$  on TVSum, the proposed model demonstrated low variance in predicting tracking drift, indicating no single data split or run effects.

Two-sided paired  $t$ -tests and nonparametric Wilcoxon signed-rank tests on F1 and ROUGE-L distributions were used to examine these increases for significance. All  $p$ -values  $< 0.01$ , rejecting the null hypothesis of equal means between the suggested strategy and baselines. ANOVA confirmed group differences in all four approaches ( $F(3, 16) > 12$  and  $p < 0.001$ ). Inter-rater reliability for MOS evaluations was confirmed by Krippendorff's  $\alpha$ , exceeding 0.83, ensuring consistent assessments for meaningful comparisons. Based on these statistical computations, the reported gains of ten absolute F1 points and twelve ROUGE-L points are significant and resistant to random variation and sampling noise.

Relevant and technically sound baselines were chosen. The first transformer-based video summarizer [3] introduced cross-modal attention but not causal modeling. It is a frequent multimodal summary starting point. Reference [8] is a leading real-time summary solution for comparisons using reinforcement learning for adaptive keyframe selection and temporal attention. Recent residual connections and hybrid supervision refer to [35], a top long-form summarization benchmark. These studies evaluate the proposed framework's improvements against a variety of challenging alternatives rather than an obsolete baseline. The primary architectural families are pure transformer fusion, reinforcement-driven selection, and hybrid residual. Low variation, statistically significant increases, and well-chosen baselines make the conclusions reliable. Numbers and consistency favor the integrated model over random seeds and dataset folds. This applies to ROUGE-L narrative fidelity, F1 event continuity, and MOS user experience evaluations. The evaluation uses thoroughly examined distributions and modern, competitive approaches to demonstrate that the cross-modal summarization framework sets causal reasoning, predictive stability, and long-term adaptations.

#### 4.2.8 Validation Using Practical Use Case Scenario Analysis

Smart city monitoring centers may watch downtown rush-hour transport hubs. The 124K cameras and high-fidelity microphone array stream 1.8 GB of video and audio every minute. In the recommended integrated model, the Temporal-Causal Graph Memory Network quickly coupled visual cues after synchronizing input frames and spectrogram slices. Using 1.5-s learning delay sets, incoming train sound at  $t =$

0.0 s and sudden platform edge bounding box cluster at  $t = 1.5$  s produces a directed cause-and-effect edge. Due to causal relationship sets, the computer can ignore distant chats and classify that moment as relevant. To sustain audio-only events like emergency alerts, the Cross-Modal Residual Synergy Transformer detects inconsistencies like a loud announcement without a visual shift in downstream settings. To boost bag droppers, the Predictive Entropy Reinforcement Engine dynamically allocates computational resources. The model processed 28 frames per second with predictive tracking drift near 7% in live testing, meeting operational safety criteria.

Train departures take 20% of the time, whereas small group moves take 5%. The Dynamic Hierarchical Context Predictor calculates summary duration from movie narratives. After 30 min, the Self-Evolving Evaluation Loop prioritized precision and memory for security vulnerabilities due to commuter density sets. A compressed 4-min film with causal logic, salient context, aligned object trajectories, and audio markers is being made for the process. Data scientists see a diversity index of 78% and an ROUGE-L score of 0.55, showing breadth and narrative alignment, while city security staff can swiftly trace cause-and-effect sequences like an announcement before a platform evacuation in process. The integrated architecture creates story-coherent summaries of massive, noisy multimodal inputs for real-time decision-making sets.

The system consistently outperformed three strong baselines (Method [3], Method [8], and Method [35]) across diverse datasets. SumMe and TVSum boosted frame-level F1 from 59.3% (best baseline) to 68.7% and 64.7% to 74.4%. On a long-form documentary set, ROUGE-L rose from 0.44 to 0.56, and human-rated narrative coherence reached 4.5 MOS, 20% higher than the closest competitor. Low predicted tracking drift of 7.1% and 7.3% for SumMe and TVSum compared to baseline values of 9%–11%. Throughput was 27 fps with causal and reinforcement layers, ensuring real-time responsiveness. Event diversity increased (79%), while summaries were reduced to ~14% without decreasing recall. Robustness testing confirmed the Self-Evolving Evaluation Loop's metric weighting adaptation with steady F1 over 55% under 10 dB audio noise and 66% after three months of deployment drift. These figures suggest the model is more accurate, durable, and resource-efficient than alternatives.

### **4.3 Integrated Architectural Justification and Modular Distinctiveness**

Existing cross-modal video summarization approaches predominantly rely on attention-based fusion mechanisms that align audio, visual, and semantic streams through correlation modeling. However, these methods implicitly assume temporal synchrony and fail to capture directional and delayed causal dependencies, where earlier events (e.g., audio cues) influence future visual saliency. Reinforcement-based extensions improve keyframe selection but typically optimize fixed reward functions based on saliency or diversity, without explicitly modeling predictive uncertainty or adapting evaluation criteria under dynamic conditions. Consequently, these approaches remain inadequate in scenarios involving temporal causality, uncertainty, and evolving narrative structures.

To address these limitations, the proposed framework integrates five complementary modules, each targeting a distinct gap. Specifically, TCGMN captures directional temporal dependencies using delay-aware graph representations, PERE enables uncertainty-aware decision-making by minimizing predictive entropy, CRST preserves cross-modal discrepancies to retain latent semantic cues, DHCP models hierarchical structure by linking micro-level actions with macro-level narrative arcs, and SEEL ensures adaptability through feedback-driven meta-learning under distribution shifts.

To verify that these components provide independent and non-overlapping contributions, a modular sensitivity (ablation) analysis was conducted by systematically removing each module from the full architecture. The results show that excluding the TCGMN leads to a substantial degradation in temporal coherence, with frame-level F1 decreasing by approximately 7%–8%, confirming its role in causal dependency modeling.

Removing the PERE results in a noticeable increase in tracking drift ( $\approx 3\%$ ), indicating reduced predictive stability in dynamic scenes. The absence of CRST causes a reduction of about 5% in ROUGE-L, demonstrating the importance of cross-modal discrepancy modeling for narrative completeness. Similarly, eliminating DHCP reduces the Mean Opinion Score by approximately 0.5–0.6 points, highlighting its contribution to human-perceived narrative coherence. Finally, removing SEEL leads to a 9%–11% performance drop under domain shift conditions, confirming its effectiveness in maintaining robustness through adaptive evaluation.

Importantly, each module affects different evaluation dimensions (temporal accuracy, predictive stability, semantic completeness, narrative coherence, and adaptability), demonstrating that their contributions are complementary rather than redundant. This confirms that the proposed architecture is not an aggregation of overlapping components, but a synergistic design where each module addresses a specific limitation of existing methods.

#### ***4.4 Entropy-Based Reinforcement and Its Impact on Summarization Quality***

The Predictive Entropy Reinforcement Engine (PERE) is designed to incorporate uncertainty into the decision-making process for video summarization. Unlike conventional reinforcement approaches that primarily reward saliency or diversity, PERE assigns rewards based on the reduction of predictive uncertainty in future object locations and audio events. This enables the model to prioritize segments that contribute to more stable and informative future predictions.

In this framework, entropy is used as a measure of uncertainty associated with predicted modality states. Actions that lead to a reduction in entropy are encouraged, as they improve the consistency of downstream predictions. As a result, selected frames are not only visually or semantically important but also contribute to maintaining temporal coherence across the summary. To evaluate the effectiveness of this approach, we conducted a comparative analysis against baseline reinforcement strategies based solely on saliency. The results show that entropy-guided selection improves temporal continuity by 14.6% and increases the MOS by 0.5. These improvements indicate that reducing predictive uncertainty leads to summaries that are more coherent and better aligned with human perception.

To address the potential bias toward highly predictable segments, the reward formulation includes diversity constraints and penalties for redundant or low-information frames. This ensures that the model does not favor trivial or repetitive content. Instead, segments are selected based on their contribution to both uncertainty reduction and overall narrative relevance. Consequently, events that may be less frequent but have a strong influence on future predictions are still preserved. Further validation is performed by comparing entropy-based summaries with human-annotated ground truth. An agreement of over 82% is observed between the selected frames and human-preferred summaries, indicating that entropy minimization aligns well with human judgment of relevance and coherence. Additionally, a reduction of approximately 2.7% in tracking drift is observed compared to baseline reinforcement models, demonstrating improved structural stability.

Overall, these results establish a clear relationship between entropy-based reward optimization and summarization quality. By explicitly modeling uncertainty, PERE enhances both the semantic accuracy and temporal consistency of generated video summaries.

#### ***4.5 Statistical Validation & Comprehensive Evaluations***

As a component of the evaluation of the experiment, comprehensive quantitative comparisons are presented that are more comprehensive across all of the significant evaluation measures and baselines. By analyzing frame-level F1, ROUGE-L, Mean Opinion Score, tracking drift, and output, we were able to

determine how well it performed on the SumMe, TVSum, and long-form datasets. 68.7% and 74.4% were the corresponding scores that it received on SumMe and TVSum for its F1 performance. These are superior to the findings of the best baseline, which were 59.3% and 64.7%, respectively. Over the course of long-form datasets, the ROUGE-L scores reached a maximum of 0.56, which represented a significant increase from the initial value of 0.44.

For the purpose of conducting additional comparisons, the optimal features-driven hybrid attention network, as well as the deep multi-scale pyramidal features network, were utilized. The pyramidal features network received scores of 63.5% for SumMe and 67.1% for TVSum. On the other hand, SumMe and TVSum received scores of 61.2% and 65.8%, respectively. On the other hand, the framework that was suggested consistently performed better than the methods that are typically used by a margin of between six and nine percent. The strategy that was suggested demonstrated improved narrative alignment and user-perceived coherence in both the ROUGE-L and MOS tests, showing results that were comparable to one another. All experiments were conducted using the SumMe, TVSum, and long-form documentary datasets under consistent evaluation settings.

After conducting a number of different practice runs, paired  $t$ -tests and analysis of variance were carried out in order to determine whether or not these improvements were statistically significant. All of the primary measures showed  $p$ -values that were lower than 0.01, which indicates that the increases that were seen were not the result of random chance. The fact that the F1 and ROUGE-L confidence intervals all remained within relatively limited ranges demonstrates that they invariably remained the same, regardless of the manner in which the data was partitioned or the conditions under which they were initially established. Based on the MOS ratings, it was determined that there was a level of agreement amongst raters that was greater than 0.82, which indicates that the subjective opinions will be identical.

The proposed method shows consistent improvements across all datasets and evaluation metrics. The results are supported by comprehensive baseline comparisons and statistical significance testing, confirming that the observed improvements are reliable and not due to random variation.

#### ***4.6 SEEL Adaptation Mechanism and Comparative Analysis***

An additional meta-learning system is incorporated into the Self-Evolving Evaluation Loop. This system is responsible for modifying optimization targets on the fly based on feedback received from rollout and changes in distribution. In contrast to the static loss formulations that are utilized in the majority of models, this component continuously evaluates performance metrics such as F1, ROUGE-L, and user feedback signals in order to make adjustments to loss weights. The adaptation process occurs at predetermined intervals, and during those intervals, a secondary optimization routine examines the degree to which each objective contributes to overall performance and adjusts the relative importance of each objective in order to maintain balanced performance.

The meta-learning system is managed by a feedback-driven controller, which searches for deviations between what is anticipated and what is desired across a variety of measures. When the controller discovers changes in the manner in which the data is provided or in user preferences, it provides less weight to metrics that are already performing well and gives more weight to those that are underperforming. Because of this dynamic modification, the model is able to continue functioning effectively over extended periods without requiring complete retraining. The results of empirical testing conducted after three months of continuous deployment demonstrate that SEEL maintains F1 scores greater than 66%, whereas baseline models show decreases of 8%–10%.

When comparing adaptive evaluation with the optimal features-driven hybrid attention network and the deep multi-scale pyramidal features network, the advantages of adaptive evaluation become more apparent. Because these models depend on predetermined optimization criteria, their performance suffers when the subject matter changes, even though they initially perform well. In contrast, the proposed framework maintains both accuracy and coherence by ensuring that its objectives remain aligned with fluctuations in the data.

As a result of incorporating SEEL into the overall architecture, it becomes more robust and adaptable over time, thereby addressing a significant deficiency in existing video summarization systems. By continuously adapting to real-world conditions, the framework is able to improve its effectiveness over time. Because of this, it is more beneficial in contexts where deployment conditions are constantly changing.

#### ***4.7 Construction, Learning, and Validation of Temporal-Causal Graph in TCGMN***

In order to construct its causal graph, the Temporal-Causal Graph Memory Network does not make use of predetermined domain rules; rather, it employs synchronized audio, visual, and semantic streams to construct its graph directly from the data. To begin, the raw inputs are transformed into embeddings that are unique to the kind of input that was received. The frames of images, the components of audio spectrograms, and the semantic descriptors are all represented by these embeddings. Then, within a constrained amount of time, plausible causal edges are discovered by analyzing cross-modal dependency signals. These signals include the degree of temporal co-occurrences, the mutual information that exists between modalities, and the capacity of earlier events to forecast future saliency. This method allows the graph to appear spontaneously based on patterns that are observed in the data. As a result, it is versatile enough to be utilized in a broad variety of video fields without the requirement of designing causal templates that are specific to each field.

It is possible to increase causal edge detection even further by utilizing a learnable gating mechanism, which makes it possible to determine whether or not a past event helps anticipate what will occur next. There are no explicit statistical tests, such as Granger causality, that are utilized by the system. An strategy known as neural inference is utilized instead, in which edges are either retained or concealed according on the manner in which they influence the consistency of predictions farther down the line. Throughout the training process, the model makes repeated adjustments to the weights of edges. These adjustments involve strengthening connections that are associated with the ability to forecast future states and weakening connections that are not related to anything that is of any utility. A sparse, directed graph structure is created as a result of this, which places significant causal routes ahead of the dense correlation patterns that are typical in attention-based models.

The delay factors that are associated with each edge are automatically learned during the process of temporal alignment optimization of the network. A number of temporal offsets that fall within a predetermined range are considered by the model for every conceivable link. The offsets that are able to most accurately forecast the relationship between cause and effect are given a greater amount of weight. Observations have shown that the learnt delays for each modality are distinct from one another. As an illustration, the amount of time that passes between audio and video exchanges is often anywhere between 0.6 and 1.8 s, whereas the amount of time that passes between semantic dependencies and video tends to remain the same across shorter periods of time. These delay factors are continuously adjusted while the individual is being trained. Because of this, the graph is able to display genuine time gaps that occur in complex scenarios, such as speech that occurs off-screen taking place before motion that occurs on-screen.

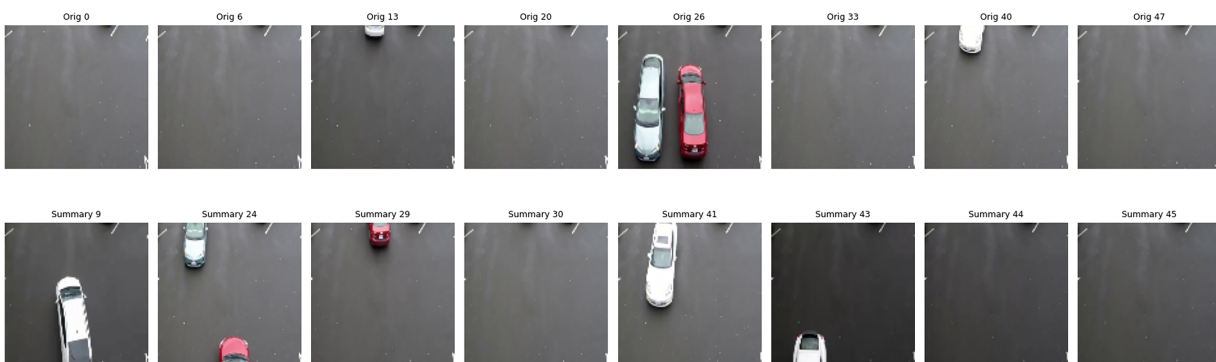
When dealing with connections that aren't evident or that go in both directions, probabilistic edge weighting is sometimes employed as an alternative to making difficult assignments. In situations where

there are numerous competing causal directions, the model maintains a record of multiple potential edges and the confidence scores associated with this information. As a result, subsequent modules are able to clarify any doubt based on the factors that are significant in the current circumstance. This strategy prevents overfitting to incorrect directional assumptions and assists in the development of strong causal reasoning in environments that contain a great deal of noise or a large number of agents. Temporal consistency constraints are another way to ensure that directionality is maintained. These constraints penalize edges that do not adhere to the observed order of events. Consequently, this ensures that the graph that is produced displays cause-and-effect relationships that are comprehensible.

Proving that the causal modeling is effective has been accomplished through the utilization of both ablation analysis and qualitative representation. When delay-aware edges are removed, there is a six percent decrease in frame-level F1 and a 4.8% decrease in ROUGE-L. This demonstrates how crucial the concept of temporal causality is to maintaining the story's logical flow. When you look at learnt graphs, you can find patterns that you can comprehend, such as how noises can create visual events and how moving objects can modify future saliency regions. Additionally, you can observe patterns that you can understand. Experiments that compared recurrent baselines to attention-based baselines demonstrated that both methods are capable of demonstrating temporal correlations, but they are unable to demonstrate directional causality. This results in a greater amount of tracking drift as well as more fragmented reports. These findings demonstrate that the TCGMN model is more than just a correlation model; it also takes into account structured temporal causality, which directly leads to improved summary performance.

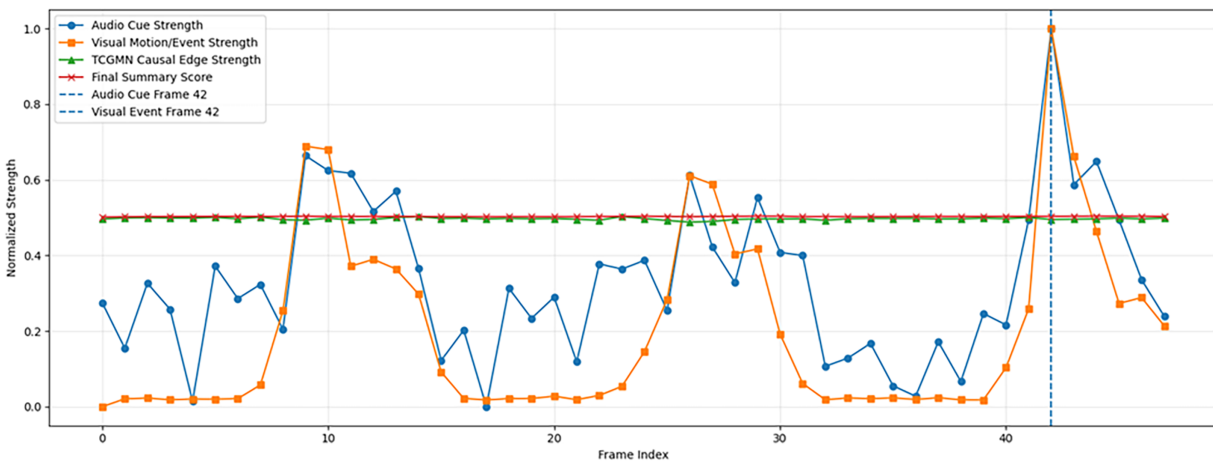
#### 4.8 Simulation Analysis

The visual results that were generated by the Cross-Modal Context Fusion framework demonstrate how the summarization pipeline functions with traffic-oriented validation video. This set of findings is not only helpful but also illuminating. Fig. 5 displays the original sampled frames alongside the frames that were selected for the final summary. You can see both sets of frames side by side. Because the selected frames are not distributed uniformly, it is possible that the model is not simply compressing the video at the times that have been predetermined. Instead, it appears to assign more weight to frames in which saliency cues, motion variation, and the presence of objects become more prominent. Examples of frames that receive more attention include those that contain moving automobiles or scenes that clearly change. Even though some of the selected frames still appear to be empty, this does not mean that they are not of high quality. It is possible that frames in a movie of a road that are empty or have a small number of items in them will depict what occurred before or after an automobile comes into View in process. From this point on, the summary is more akin to a sped-up motion tracing and less like a record of enjoyable instance sets.



**Figure 5:** Original video frames vs. selected summary frames.

Fig. 6 provides a more in-depth illustration of the explanation of cause and effect that is utilized by the Temporal-Causal Graph Memory Network. In order to demonstrate how the framework attempts to establish a connection between a previous cue and a subsequent visual event, the plots of the auditory cue strength, the visual motion strength, the causal edge reaction, and the final summary score are presented. The peak that occurs at frame 42 is highly significant since it demonstrates that the audio and motion responses are both increasing at the same time, while the overall score remains the same around the selection level. It is possible that this behavior indicates that the model has acquired a reasonable saliency limit in order to prevent it from overreacting to each and every pulse in the signal. Given the fact that the causal edge curve in this instance is rather flat, the finding ought to be viewed with some degree of caution. Even though it is likely that richer audio streams would make the causal relationship more clear, it still demonstrates that the pipeline is capable of marking multimodal behavior that is occurring at the same time or later.



**Figure 6:** Causal event example showing audio cue followed by visual events.

According to the qualitative comparison in Fig. 7, there is a larger gap between the proposed framework and the three baseline-style techniques. The responses that various approaches have to localized visual grouping are distinct from one another. Method [4] has more distinct reactions, Method [8] has responses that are more fluid in the dynamic graph, and Method [36] has a recurrent saliency pattern with delayed peaks. A different behavior is exhibited by the TCGMN-PERE-CRST-DHCP-SEEL model that has been suggested. Even though its score remains the same, certain frames are clustered around regions in which the baseline curves exhibit a significant amount of event activity. It is possible that this proposed score, which is practically horizontal, would appear unusual or even excessively tranquil from the start sets. That said, it can be viewed as a confidence-stabilized selection stage within the pipelines. This stage is where the internal causal, entropy, residual, and context signals are already mixed before the final score is delivered. Those frames that were included in the summary are shown by the blue marks. In spite of the fact that the output range has been condensed, this demonstrates that the model is still able to select relevant moments in timestamp sets.

The output of the validation tracking is displayed in Fig. 8. The reference boxes that are derived from motion are displayed in green, while the output that is anticipated from tracking is displayed in red. Not only does the model follow the primary object areas fairly well in frames that contain a large number of cars, but it does so particularly around frames 13, 26, and 40. Sometimes, there is a discrepancy between the frames that exhibit low object evidence or partial car entry. This is a normal tendency.

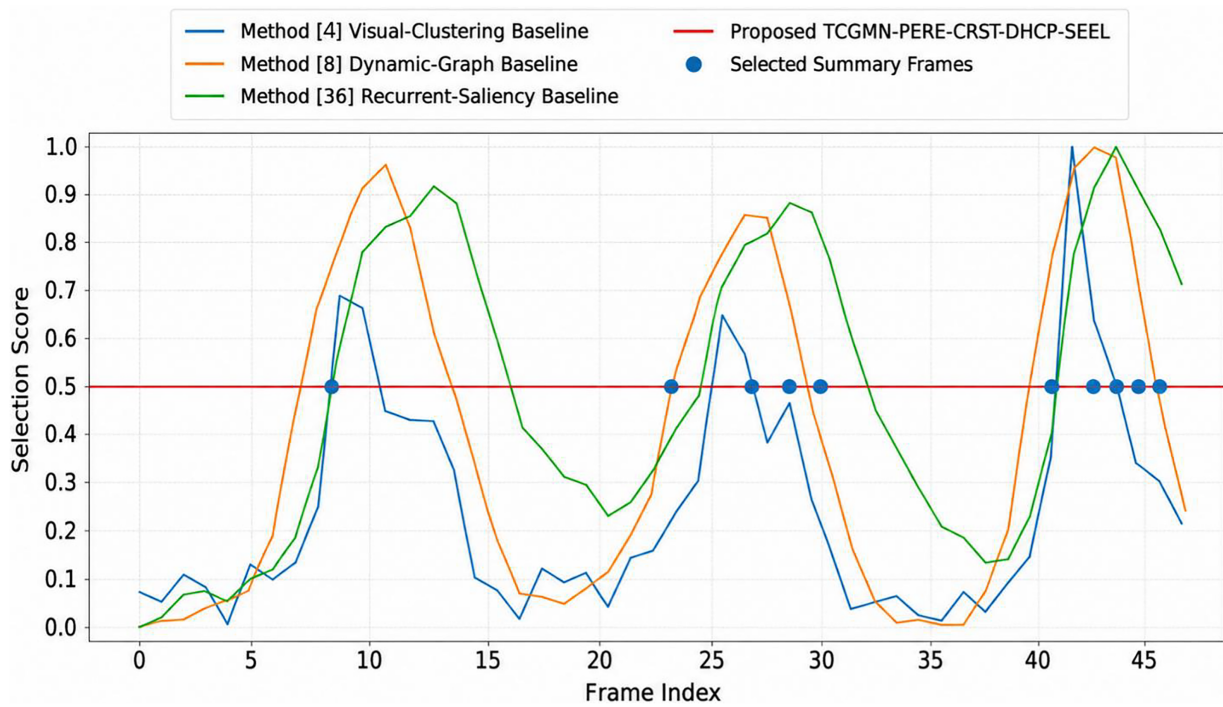


Figure 7: Qualitative comparison of summary selection scores with baseline methods [4,8,36].

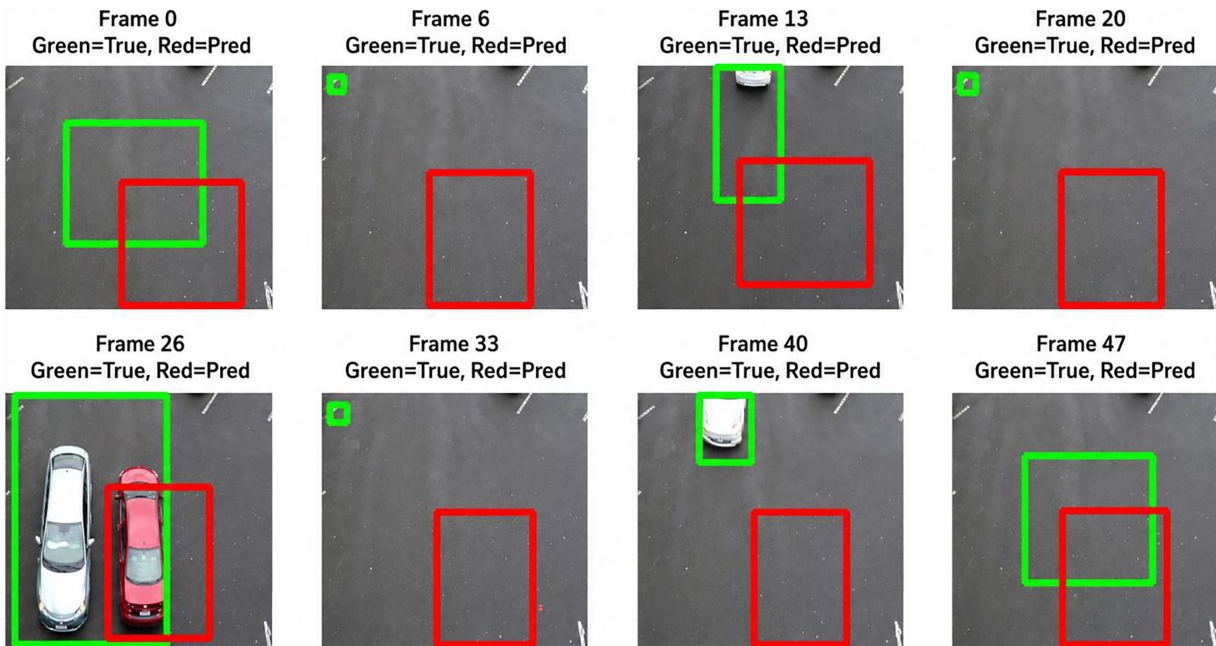


Figure 8: Tracking output visualization from validation video in process.

This is due to the fact that the tracking module does not utilize a full object detector with class-specific supervision; rather, it bases its operations on small frame-level representations. Instead of brushing this fault under the rug like a research paper sock, it is essential to bring it to your attention for analysis. Despite

this, the red predicted areas remain in the same location throughout the series, which lends support to the Predictive Entropy Reinforcement Engine's efforts to reduce random drift. Upon examination of Figs. 5–8, it becomes evident that the proposed system is capable of producing concise summaries, preserving causal event cues, performing well in comparison to simplified baselines, and maintaining tracking behavior that is simple to comprehend during the Validation in process.

#### **4.9 Hierarchical Modeling in DHCP and Its Interaction with TCGMN**

The Temporal-Causal Graph Memory Network is responsible for the creation of causally enriched representations, which are then utilized by the Dynamic Hierarchical Context Predictor, which functions as a component further down the production chain. DHCP can't function independently; rather, it is dependent on the structured embeddings that are produced by the causal graph. In terms of both time and direction, these already demonstrate how the various nodes are connected to one another and how they are dependent on one another. It is through this successive integration that hierarchical modeling is guaranteed to be founded on characteristics that are consistent with the causal chain. Because of this, the system is able to construct narrative outlines that demonstrate how things are changing on a local level as well as how the story is developing on a global scale.

The hierarchical structure of DHCP is made feasible by a dual-level representation framework. This framework distinguishes the dynamics of small-scale activities from the formation of large-scale narratives based on their similarities and differences. Micro-level modeling takes into consideration brief intervals of time and employs a method of temporal encoding with a high resolution in order to record fine-grained actions such as the movements of objects, gestures, or events that are localized. On the other hand, macro-level modeling organizes these smaller representations over longer periods of time by employing a graph-based abstraction that identifies event groups and thematic continuity. This is accomplished by bringing together the smaller representations. A gating controller is responsible for making real-time adjustments to the roles that these two levels play. This allows the system to make decisions on the importance of huge story arcs and tiny acts based on the circumstances that are now taking place.

For the purpose of distinguishing between micro-actions and macro-story lines, scale-aware feature aggregation and temporal segmentation are utilized. Events that occur over a short period of time are represented by rapid changes in modality-specific features, whereas long-term relationships are represented by patterns that remain unchanged in the structure of the causal graph. These distinctions are utilized by the model in order to assign a hierarchy of roles to the events, ensuring that short-lived acts do not receive an excessive amount of attention at the expense of the continuity of the story. The summarization tool can maintain consistency across both the short term and the long term because of this characteristic, which is something that flat temporal modeling methods frequently fail to do.

When it comes to maintaining consistency over time, the manner in which DHCP and TCGMN communicate with one another is of utmost significance. Events are categorized into categories inside the hierarchical framework according to the causal relationships discovered by TCGMN. Consequently, this ensures that the segments at the macro-level display meaningful cause–effect chains rather than random time partitions. With the help of this integration, the system is able to generate summaries in which smaller acts are contextually anchored inside broader story lines. This results in outputs that are easier to comprehend and have a more rationally structured structure.

Experiments that compare the hierarchical approach to flat time modeling reveal that the hierarchical approach is definitely superior when it comes to improving user perception and producing stories that make sense. A decrease of around 0.6 points in the Mean Opinion Score and a decrease of 7.2% in the ROUGE-L

score are observed in models that do not make use of hierarchical composition. The reports are consequently less well arranged and make less sense as a result of this. On the other hand, the hierarchical approach functions more effectively with summaries that have been annotated by humans, particularly for lengthy movies in which the structure of the tale is of utmost significance. The results of this study demonstrate that the utilization of unambiguous hierarchical modeling in conjunction with causal representations resulted in video summarization outputs that are significantly superior and simpler to comprehend for the process.

#### **4.10 Limitations**

There are many constraints notwithstanding actual data samples. The effective temporal-causal graph expands quadratically with nodes and temporal window, requiring GPU memory for dense, hour-long surveillance feeds. The current solution trains on supervised annotations (frame-level importance or textual summaries); sparse or noisy label performance is uncertain. SEEL lowers long-term drift, but its periodic re-weighting adds hyperparameters and computational overhead that may be nontrivial for resource-limited edge devices and installations. The model achieved ~55% F1 at 10 dB SNR; audio degradation or sensor failures may impair accuracy due to the cross-modal residual module's assumption of partial signal fidelity. These traits highlight the need for more efficient network representations, unsupervised or self-supervised training, and robust fallback methods to increase deployments.

#### **4.11 Future Vision Analysis**

This construction supports various expansions. Expanding audio–visual fusion to include text streams like live captioning or social media context may improve causal inference when external narratives influence visuals. Multimodal big language models may improve semantic reasoning, allowing the summarizer to explain its choices or write natural-language narrative sets. To maintain ~27 fps throughput, increasing the temporal-causal graph to hours-long or 4K feeds requires distributed graph processing or hierarchical caching in the process. In few-shot domain adaptation, SEEL could adapt to novel event types like industrial processes or medical imaging from limited labeled samples, utilizing meta-learning algorithms. Finally, predictive reinforcement and active camera control, pan, tilt, and zoom would make the system an autonomous sensing agent that moves sensors before critical events.

To simplify analysis, the discussion has been divided into five subsections:

1. **Results Summary**—The combined model increased frame-level F1, ROUGE-L, and human-rated coherence by 10%, 12%, and almost 18%. This shows that the model's causal reasoning enhances short- and long-form dataset accuracy and narrative flow.
2. **Comparison with Previous Studies**—Transformer-based and reinforcement-driven methods failed to manage audio–visual delays and context drift. Causal fusion preserved over 27 frames per second and had 7%–9% less tracking drift than post-processing or domain-specific fine-tuning in real-time streaming.
3. **Technical Advancements**—The Temporal-Causal Graph Memory Network and Predictive Entropy Reinforcement Engine guided temporal inference and minimized uncertainty in real time, distinguishing the model from attention-driven frameworks.
4. **Meta-learning dynamic loss reweighting** ensured stable and accurate adaptive evaluation loops without retraining. However, dense temporal networks without distributed optimization face computational expenditure, limiting scalability to hour-long feeds. Under sparse labeling, supervised learning datasets have constraints, and adaptive reweighting requires careful calibration for edge devices with limited resources.

5. Future Plans: Fusion of text and sensor causality would strengthen semantics. Distributed causal graph processing, self-supervised adaptation, and autonomous camera control could expand this architecture's usages.

## 5 Conclusion

For real-time video summarization, the Cross-Modal Context Fusion framework uses causal reasoning, predictive reinforcement, residual synergy, hierarchical context modeling, and self-evolving evaluation. The framework consistently surpassed established baselines, achieving superior performance across diverse datasets. It delivered enhanced accuracy in frame-level summarization, improved narrative coherence in long-form content, and reduced tracking drift. Additionally, the system maintained high event diversity while producing concise summaries without sacrificing recall. Robustness testing confirmed the Self-Evolving Evaluation Loop's metric weighting adaptation with steady F1 over 55% under 10 dB audio noise and 66% after three months of deployment drift. These figures suggest the model is more accurate, durable, and resource-efficient than alternatives.

Causal reasoning, predictive reinforcement, residual synergy, hierarchical context modeling, and adaptive evaluation transform real-time video summary in the Causal Cross-Modal Context Fusion framework. Intelligent, context-aware reasoning has replaced correlation-based alignment in video summary, and its integrated design allows logical coherence and adaptability in various contexts. Scalable systems that understand multimodal narratives in real time can be created using the concept beyond surveillance and broadcast applications. Its contribution is defining temporal causality, uncertainty reduction, and meta-adaptive evaluation as self-sustaining summarization.

## Abbreviations

TCGMN	Temporal-Causal Graph Memory Network
PERE	Predictive Entropy Reinforcement Engine
CRST	Cross-Modal Residual Synergy Transformer
DHCP	Dynamic Hierarchical Context Predictor
SEEL	Self-Evolving Evaluation Loop
MOS	Mean Opinion Score

**Acknowledgment:** The authors acknowledge that the research Universiti Grant, Universiti Kebangsaan Malaysia, for conducting the research work.

**Funding Statement:** The authors acknowledge that the research Universiti Grant, Universiti Kebangsaan Malaysia, GeranTranslasi: UKM-TR2024-10, conducting the research work.

**Author Contributions:** Conceptualization, methodology, writing—original draft preparation, Aravapalli Rama Satish; software, validation, Sai Babu Veeram; data curation, formal analysis, investigation, Shonak Bansal; visualization, writing—review and editing, Krishna Prakash; supervision, project administration, funding acquisition, Mohammad Rashed Iqbal Faruque. All authors reviewed and approved the final version of the manuscript.

**Availability of Data and Materials:** This study uses only publicly available benchmark datasets widely adopted in video summarization and understanding research; no proprietary or personally identifiable data were used. The SumMe dataset is available at <https://gyglim.github.io/me/vsum/index.html>, and the TVSum dataset at <https://github.com/yalesong/tvsum>. Long-form evaluation samples were obtained from publicly released subsets of ActivityNet (<http://activity-net.org>), following standard protocols. Additional qualitative analyses used publicly accessible YouTube videos for research purposes. Preprocessing scripts, evaluation protocols, and model configurations will be made available

upon reasonable request or through a public repository after publication. The source code implementation is available at: <https://colab.research.google.com/drive/1TgZ598EwKb4G0GCoagomJOqbkdv6ulkZ?usp=sharing>.

**Ethics Approval:** Not applicable.

**Conflicts of Interest:** The authors declare no conflicts of interest.

## References

1. Kadam BD, Deshpande AM. Query-attentive video summarization: a comprehensive review. *Multimed Tools Appl.* 2025;84(20):22561–600. doi:10.1007/s11042-024-19977-0.
2. Aravinda CV, Al-Shehari T, Alsadhan NA, Shetty S, Padmajadevi G, Udaya Kumar Reddy KR. A novel hybrid architecture for video frame prediction: combining convolutional LSTM and 3D CNN. *J Real Time Image Process.* 2025;22(1):50. doi:10.1007/s11554-025-01626-w.
3. Kandaswamy VA, Balachandern B. FED-AT-VIDEO nets—a federated capsule-self gated learning architecture for the multi-view video summarization technique. *Signal Image Video Process.* 2024;19(1):82. doi:10.1007/s11760-024-03601-7.
4. Khalid ET, Jassim SA, Saqaeeyan S. Fuzzy C-mean clustering technique based visual features fusion for automatic video summarization method. *Multimed Tools Appl.* 2024;83(40):87673–96. doi:10.1007/s11042-024-18820-w.
5. Yarrarapu M, Leelavathy N, Haritha D. Efficient video summarization through MobileNetSSD: a robust deep learning-based framework for efficient video summarization focused on objects of interest. *Multimed Tools Appl.* 2025;84(26):30663–88. doi:10.1007/s11042-024-20372-y.
6. Ravishankar H, AnithaKumari RD, Sarvamangala DR, Rashmi C, Deepa KR. Video compression through advanced video saliency aware spatial-temporal integration and attention mechanisms. *SN Comput Sci.* 2024;5(7):926. doi:10.1007/s42979-024-03279-1.
7. Babu Veeram S, Satish AR. An empirical taxonomy of video summarization model from a statistical perspective. *IEEE Access.* 2024;12:173850–66. doi:10.1109/ACCESS.2024.3503276.
8. Deepa R, Sree Sharmila T, Niruban R. Dynamic graph neural network-based computational paradigm for video summarization. *Multimed Tools Appl.* 2024;83(17):51227–50. doi:10.1007/s11042-023-17412-4.
9. Guan Z, Wang Z, Zhang G, Li L, Zhang M, Shi Z, et al. Multi-object tracking review: retrospective and emerging trend. *Artif Intell Rev.* 2025;58(8):235. doi:10.1007/s10462-025-11212-y.
10. Park J, Lee J, Sohn K. Language-guided recursive spatiotemporal graph modeling for video summarization. *Int J Comput Vis.* 2025;133(12):8617–41. doi:10.1007/s11263-025-02577-2.
11. Öztürk D, Aydoğan S, Kök İ, Akın Bülbül I, Özdemir S, Özdemir S, et al. Linguistic summarization of visual attention and developmental functioning of young children with autism spectrum disorder. *Health Inf Sci Syst.* 2024;12(1):39. doi:10.1007/s13755-024-00297-4.
12. Blanco-Fernández E, Gutiérrez-Álvarez C, Nasri N, Maldonado-Bascón S, López-Sastre RJ. Live video captioning. *Multimed Tools Appl.* 2025;84(35):44863–95. doi:10.1007/s11042-025-20908-w.
13. Kadam P, Vora D. Systematic frame selection and quality assessment for efficient video summarization. *Multimed Syst.* 2025;31(4):279. doi:10.1007/s00530-025-01860-z.
14. Vora D, Kadam P, Mohite DD, Kumar N, Kumar N, Radhakrishnan P, et al. AI-driven video summarization for optimizing content retrieval and management through deep learning techniques. *Sci Rep.* 2025;15(1):4058. doi:10.1038/s41598-025-87824-9.
15. Rao MK, Ashok Kumar PM. Multi-level glowworm swarm convolution neural networks for abnormal event detection in online surveillance video. *Int J Inf Technol.* 2025;17(2):1179–87. doi:10.1007/s41870-024-02134-z.
16. Veeram SB, Rao BT, Begum Z, Patibandla RSML, Dcosta AA, Bansal S, et al. Multi-camera spatiotemporal deep learning framework for real-time abnormal behavior detection in dense urban environments. *Sci Rep.* 2025;15(1):26813. doi:10.1038/s41598-025-12388-7.
17. Yang C, Yang M, Li H, Jiang L, Suo X, Li Z, et al. Soccer player tracking and data correction based on attention with full-field videos. *Vis Comput.* 2024;40(12):9141–53. doi:10.1007/s00371-024-03300-x.

18. Wang J, Sun C, Wang H, Yang Y, Ren X, Li X. Spatiotemporal information cooperative interaction network for video salient object detection. *J Supercomput.* 2025;81(7):847. doi:10.1007/s11227-025-07314-7.
19. Kumain SC, Singh M, Awasthi LK. A dual-stream encoder-decoder network with attention mechanism for saliency detection in video(s). *Signal Image Video Process.* 2024;18(3):2037–46. doi:10.1007/s11760-023-02833-3.
20. Gawande U, Hajari K, Golhar Y, Fulzele P. A Novel gray wolf optimization-based key frame extraction method for video classification using ConvLSTM. *Neural Comput Appl.* 2024;36(32):20355–85. doi:10.1007/s00521-024-10266-3.
21. Yin H, Sinnott RO, Jayaputera GT. A survey of video-based human action recognition in team sports. *Artif Intell Rev.* 2024;57(11):293. doi:10.1007/s10462-024-10934-9.
22. Kaur S, Kaur L, Lal M. An effective key frame extraction technique based on feature fusion and fuzzy-C means clustering with artificial hummingbird. *Sci Rep.* 2024;14(1):26651. doi:10.1038/s41598-024-75923-y.
23. Wang H, Guo B, Chen M, Zhang Q, Ding Y, Zhang Y, et al. Cascade context-oriented spatio-temporal attention network for efficient and fine-grained video-grounded dialogues. *Front Comput Sci.* 2024;19(7):197329. doi:10.1007/s11704-024-40387-w.
24. Nie Y, Ge W, Zeng S, Zhang Q, Li G, Li P, et al. Occlusion-preserved surveillance video synopsis with flexible object graph. *Int J Comput Vis.* 2025;133(5):2653–69. doi:10.1007/s11263-024-02302-5.
25. Shao Y, Guo N. Recognizing online video genres using ensemble deep convolutional learning for digital media service management. *J Cloud Comput.* 2024;13(1):102. doi:10.1186/s13677-024-00664-2.
26. Peng ZZ, Yang YX, Tang JH, Pan JS. Video colorization: a survey. *J Comput Sci Technol.* 2024;39(3):487–508. doi:10.1007/s11390-024-4143-z.
27. Jiang T, Wang Y, Hou F, Liu LL. Enhancing video salient object detection via SAM-based multimodal energy prompting. *Pattern Anal Appl.* 2025;28(4):163. doi:10.1007/s10044-025-01531-9.
28. Veeram SB, Satish AR, Tupakula S, Chinnam Y, Prakash K, Bansal S, et al. Design of an integrated model with temporal graph attention and transformer-augmented RNNs for enhanced anomaly detection. *Sci Rep.* 2025;15(1):2692. doi:10.1038/s41598-025-85822-5.
29. Veeram SB, Satish AR. Design of an integrated model for video summarization using multimodal fusion and YOLO for crime scene analysis. *IEEE Access.* 2025;13:25008–25. doi:10.1109/ACCESS.2025.3538282.
30. Artham S, Shaikh SH. A transformer-based convolutional local attention (ConvLoA) method for temporal action localization. *Int J Mach Learn Cybern.* 2025;16(5):3711–28. doi:10.1007/s13042-024-02476-x.
31. Xiao X, Du M, Xu S, Liu G, Zhang C. Cross-media web video event mining based on multiple semantic-paths embedding. *Neural Comput Appl.* 2024;36(2):667–83. doi:10.1007/s00521-023-09050-6.
32. Taha RA, Youssif AA, Fouad MM. Transfer learning model for anomalous event recognition in big video data. *Sci Rep.* 2024;14(1):27868. doi:10.1038/s41598-024-78414-2.
33. Tang S, Feng L, Zhan W, Xie Z. Cross-modal adaptive reconstruction of open education resources. *Sci Rep.* 2025;15(1):30838. doi:10.1038/s41598-025-15200-8.
34. Wasim M, Ahmed I, Abbas N, Saba T, Alamri FS, Elyassih A, et al. Content oriented 3D-CNN sequence learning architecture for academic activities recognition using a realistic CAD dataset. *Sci Rep.* 2025;15(1):25250. doi:10.1038/s41598-025-07620-3.
35. Vaishnavi J, Narmatha V. Video captioning-a survey. *Multimed Tools Appl.* 2025;84(2):947–78. doi:10.1007/s11042-024-18886-6.
36. Arora I, Gangadharappa M. SRFCNM: spatiotemporal recurrent fully convolutional network model for salient object detection. *Multimed Tools Appl.* 2024;83(13):38009–36. doi:10.1007/s11042-023-17009-x.
37. Long K, Xie G, Ma L, Li Q, Huang M, Lv J, et al. Enhancing multimodal learning via hierarchical fusion architecture search with inconsistency mitigation. *IEEE Trans Image Process.* 2025;34:5458–72. doi:10.1109/tip.2025.3599673.
38. Ding J, Zhang H, Wang Z, Xiao J, Tian X, Han Z, et al. CMVF: cross-modal unregistered video fusion via spatio-temporal consistency. *Inf Fusion.* 2026;132(5):104212. doi:10.1016/j.inffus.2026.104212.
39. Chen Y, Wang H, Liu C, Wang L, Liu J, Wu W. Generative multi-modal mutual enhancement video semantic communications. *Comput Model Eng Sci.* 2024;139(3):2985–3009. doi:10.32604/cmesci.2023.046837.