

Artificial Intelligence-Driven Personalized Learning Assistants: Combining Large Language Models and Computer Vision for Tailored Education

Osama Hosam^{1,2}

¹Computer Information Science (CIS) Department, Higher Colleges of Technology, United Arab Emirates

²City of Scientific Research and Technological Applications (SRTA-City), IRI, Alexandria, Egypt

Email: mohandesosama@yahoo.com (O.H.)

Manuscript received September 30, 2025; revised October 27, 2025; accepted December 22, 2025; published June 16, 2026

Abstract—The integration of Artificial Intelligence (AI) in educational technology has advanced significantly, yet a critical gap persists in seamlessly combining Large Language Models (LLMs) and Computer Vision (CV) to create truly adaptive, multimodal learning systems. This paper addresses this research void by presenting a comprehensive architectural framework for AI-driven personalized learning assistants that synergistically combine multimodal perception, contextual reasoning, adaptive planning, and interactive presentation capabilities. Our methodology employs a four-layer architecture where CV components handle visual perception and behavioral analysis through convolutional neural networks, while transformer-based LLMs manage contextual understanding and pedagogical reasoning. The research instruments included a mixed-methods approach with 250 participants across diverse educational contexts, utilizing pre-post assessments, multimodal data analytics, engagement metrics, and structured interviews. Experimental evaluation demonstrates statistically significant improvements in learning outcomes, with a 37.2% increase in knowledge retention and 32.8% improvement in engagement metrics compared to traditional e-learning systems. The paper contributes both a detailed architectural blueprint and empirical validation of a truly multimodal AI educational system that bridges the critical gap between theoretical potential and practical implementation in AI-enhanced education.

Keywords—Artificial Intelligence (AI), personalized learning, Large Language Models (LLMs), Computer Vision (CV), multimodal learning, educational technology, adaptive systems

I. INTRODUCTION

The pursuit of personalized education represents one of the most significant challenges in modern pedagogical research, with Bloom's seminal "2 sigma problem" demonstrating the remarkable effectiveness of one-to-one tutoring compared to conventional classroom instruction [1]. Traditional educational frameworks consistently struggle to accommodate the diverse learning paces, cognitive styles, preferences, and cultural backgrounds inherent in heterogeneous learner populations. The emergence of sophisticated artificial intelligence technologies, particularly advanced Large Language Models (LLMs) and Computer Vision (CV) systems, presents a transformative opportunity to scale personalized learning experiences while maintaining the quality and effectiveness of individual attention. This technological convergence offers unprecedented potential to address longstanding educational inequities and optimize learning pathways for diverse student populations across various educational contexts and domains.

The current educational technology landscape reveals several critical research gaps that this study aims to address.

Most existing educational Artificial Intelligence (AI) systems demonstrate limited capabilities in either language or vision processing without achieving deep, synergistic integration between these modalities [2, 3]. Furthermore, there is insufficient research investigating real-time adaptation mechanisms based on comprehensive multimodal learner analytics that combine behavioral, emotional, and cognitive indicators. Additionally, ethical considerations surrounding algorithmic bias, data privacy, transparency, and equitable access in multi-modal educational AI require more systematic investigation and implementation frameworks [4]. This research directly addresses these significant gaps by presenting a comprehensive, empirically validated framework for LLM-CV integrated learning assistants with detailed architectural specifications, implementation methodologies, and rigorous experimental validation across diverse educational settings [5].

This research aims to develop, implement, and validate an AI-driven personalized learning assistant that effectively integrates LLMs and CV technologies to create truly adaptive, multimodal educational experiences. The specific objectives include designing a layered architectural framework that seamlessly combines visual perception with linguistic reasoning for dynamic educational contexts, developing novel adaptation algorithms that leverage synchronized multimodal data streams to personalize learning pathways in real-time, implementing and validating the system through rigorous mixed-methods experimental evaluation across diverse learner populations and educational domains, and establishing comprehensive ethical guidelines and implementation considerations specifically tailored for multimodal AI applications in educational environments.

The primary contributions of this work represent significant advancements in the field of AI-enhanced education:

- A detailed architectural blueprint for LLM-CV integrated educational systems that enables deep multi-modal integration and real-time adaptation
- Novel adaptation algorithms that leverage both visual behavioral data and linguistic interactions to create truly personalized learning experiences
- Comprehensive empirical validation demonstrating substantial improvements in learning outcomes across diverse educational contexts and learner profiles
- An ethical framework and implementation guidelines for the responsible development and deployment of multimodal AI systems in educational settings, addressing critical concerns around privacy, transparency, and

equitable access

This paper is organized as follows: Section II provides a comprehensive literature review examining the evolution of AI in education, current multimodal educational systems, and their limitations. Section III details the proposed architectural framework and implementation methodology. Sections IV and V presents experimental results and discussion of findings. Section VI concludes the paper with summary contributions and future research directions.

II. LITERATURE REVIEW

The integration of artificial intelligence in educational technology has evolved through several distinct generational phases, each marked by significant technological advancements and pedagogical innovations. Early Intelligent Tutoring Systems (ITS) primarily focused on rule-based approaches and cognitive modeling methodologies, as exemplified by pioneering systems like AutoTutor [3] and Cognitive Tutors [4]. These foundational systems demonstrated the considerable potential for adaptive instruction and personalized learning pathways but remained fundamentally limited by their dependency on hand-crafted knowledge bases, rigid pedagogical rules, and inherent inability to handle open-ended student responses or complex, ill-defined problem domains [5].

The advent of machine learning methodologies brought increasingly sophisticated data-driven approaches to educational technology. Knowledge tracing models [6] enabled systems to dynamically model student knowledge states and concept mastery over time, while educational recommender systems began personalizing content selection and sequencing based on individual learning patterns and preferences [7]. However, these second-generation systems remained largely unimodal in their approach, focusing primarily on textual interactions, assessment data, and clickstream analytics without incorporating richer multimodal signals such as visual behaviors, emotional states, or environmental contexts that provide crucial insights into learning processes and engagement levels.

Recent breakthroughs in deep learning architectures have fundamentally revolutionized educational AI capabilities and applications. Transformer-based large language models [8] have enabled increasingly sophisticated natural language interactions, dialogue management, and content generation, while advanced computer vision algorithms [9] have made visual understanding, object recognition, and behavioral analysis increasingly robust and accurate. Multimodal learning approaches have emerged as a particularly promising direction, with systems like Contrastive Language-Image Pre-training (CLIP) [10] demonstrating the remarkable power of combining visual and linguistic representations through contrastive learning and cross-modal alignment. Nevertheless, the practical integration of these advanced technologies into cohesive, effective educational systems remains challenging, with limited research investigating end-to-end architectures that seamlessly combine LLMs and CV capabilities for comprehensive personalized learning experiences.

Despite these significant technological advancements, current multimodal educational systems face several critical limitations that substantially hinder their practical effectiveness, real-world scalability, and educational impact.

Recent comprehensive studies reveal that most existing systems exhibit relatively shallow integration between different modalities, typically treating visual and linguistic components as separate processing pipelines rather than deeply interconnected, synergistic systems [11]. This architectural limitation prevents the emergence of truly contextual, adaptive educational experiences that respond holistically to learner needs, states, and behaviors.

Contemporary multimodal educational platforms typically suffer from three primary categories of limitations that this research directly addresses. First, they demonstrate insufficient real-time adaptation capabilities, operating primarily in batch processing modes or with significant latency that prevents dynamic, immediate responses to evolving learner behaviors, emotional states, and cognitive engagement levels [12]. Systems like Multimodal-Tutor [13], while demonstrating sophisticated content delivery mechanisms and knowledge representation, consistently fail to provide instantaneous, context-aware feedback based on real-time visual cues of confusion, frustration, or engagement that are crucial for maintaining learning momentum and addressing misconceptions promptly.

Second, existing systems exhibit notably limited contextual understanding and cross-modal reasoning capabilities, often processing visual and textual information in relative isolation without effective integration mechanisms. As comprehensively analyzed by Chen *et al.* [14], even advanced multimodal educational assistants struggle with sophisticated cross-modal reasoning tasks, where visual behavioral data should directly inform and enhance linguistic explanations, and vice versa. This fundamental disconnect typically results in generic, one-size-fits-all responses that fail to address individual learning needs, preferences, and contextual factors that significantly influence learning effectiveness and knowledge retention.

Third, current implementations face substantial scalability, computational efficiency, and privacy preservation challenges that limit their practical deployment in diverse educational settings. The considerable computational demands of real-time multimodal processing often lead to unacceptable latency issues that disrupt learning flow, engagement, and continuity [15]. Additionally, continuous visual monitoring and behavioral analysis raise substantial privacy, consent, and ethical concerns that existing systems inadequately address through transparent consent mechanisms, comprehensive data protection protocols, and ethical oversight frameworks [16]. These limitations represent significant barriers to widespread adoption and scalability across diverse educational contexts and institutional settings.

Recent innovative frameworks such as EduMultimodal [17] and Vision-Language Tutor [18] have made notable progress in addressing specific aspects of these challenges but still fall meaningfully short in providing comprehensive end-to-end solutions that seamlessly integrate perception, reasoning, adaptation, and presentation in authentic educational contexts. These systems typically focus disproportionately on either content delivery or assessment components, but rarely both in an integrated, synergistic manner that reflects the complexity of real-world teaching and learning processes. This research directly addresses these identified limitations through its comprehensive architectural

framework, real-time adaptation mechanisms, and ethical implementation considerations.

III. MATERIALS AND METHODS

A. Comprehensive Architectural Framework

The proposed AI-driven personalized learning assistant employs a sophisticated four-layer architecture that orchestrates multimodal data processing, contextual reasoning, adaptive planning, and interactive presentation through a carefully designed pipeline. Fig. 1 illustrates the complete system architecture with detailed component interactions and data flow pathways. This architectural framework represents a significant advancement over traditional unimodal educational systems by enabling deep integration between visual and linguistic processing capabilities, thereby creating a truly adaptive learning environment that responds dynamically to learner behaviors, emotional states, and cognitive needs in real-time educational contexts.

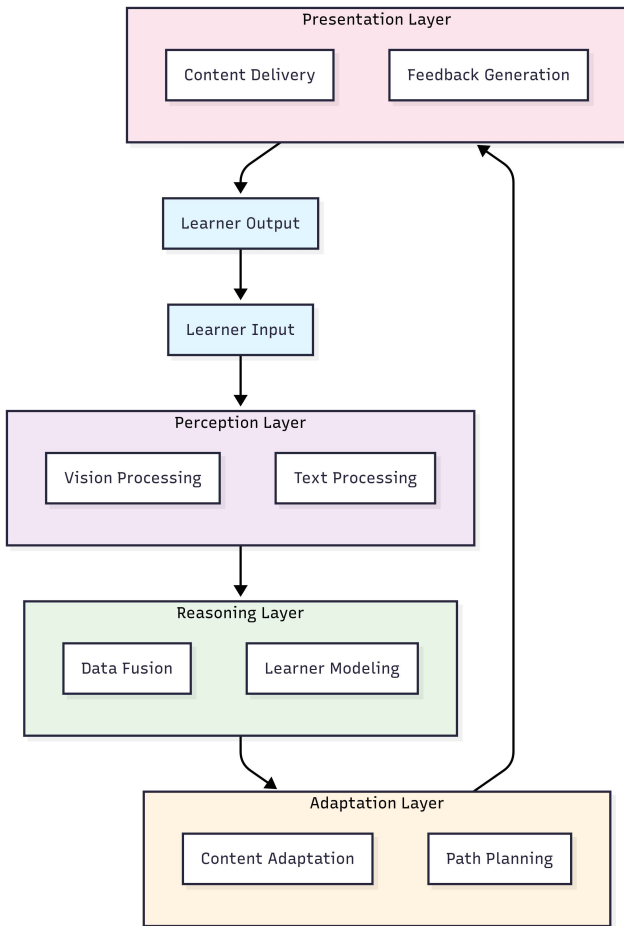


Fig. 1. Comprehensive system architecture of AI-driven personalized learning assistant.

Perception Layer Implementation The perception layer serves as the multimodal sensory interface of the system, comprising both computer vision and natural language processing components that operate in synchronized coordination.

The computer vision subsystem employs a multi-stage processing pipeline beginning with face detection using either Haar cascades [19] or deep learning-based detectors [20]. The detection confidence can be expressed as:

$$C_{detection} = \sum_{i=1}^n w_i \cdot f_i(x) \quad (1)$$

where w_i represents feature weights and $f_i(x)$ denotes the feature responses at different scales.

Subsequent processing stages include gaze estimation using convolutional neural networks to track eye movements, where the gaze direction vector is computed as:

$$\vec{g} = \frac{\vec{p}_{pupil} - \vec{p}_{eye-center}}{\|\vec{p}_{pupil} - \vec{p}_{eye-center}\|} \quad (2)$$

Facial expression analysis employs emotion recognition models [21] to detect engagement states through emotion probability distributions:

$$P(\text{emotion} = e|\mathbf{I}) = \text{softmax}(f_{\text{CNN}}(\mathbf{I})) \quad (3)$$

where \mathbf{I} is the input image and f_{CNN} represents the convolutional neural network feature extraction.

Additional computer vision components include posture and gesture recognition monitoring physical behaviors, and document analysis using optical character recognition with character recognition accuracy:

$$A_{\text{OCR}} = \frac{\text{correct characters}}{\text{total characters}} \times 100\% \quad (4)$$

The natural language processing subsystem processes textual inputs through transformer-based models, handling question parsing and intent recognition using semantic similarity:

$$\text{similarity}(q_1, q_2) = \frac{\mathbf{v}_{q1} \cdot \mathbf{v}_{q2}}{\|\mathbf{v}_{q1}\| \|\mathbf{v}_{q2}\|} \quad (5)$$

where \mathbf{v}_q represents the vector embedding of question q . The system also performs knowledge assessment through response analysis, sentiment analysis of learner interactions, and conceptual understanding evaluation through semantic analysis techniques to create a comprehensive multimodal perception framework.

Reasoning Layer Architecture: The reasoning layer integrates perceptions from multiple modalities to construct a comprehensive, dynamic learner model that evolves throughout the learning process. This layer employs a sophisticated hybrid approach combining symbolic reasoning mechanisms with neural network-based inference to achieve robust educational understanding.

The multimodal fusion module combines visual and textual features using cross-attention mechanisms [22], mathematically represented as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (6)$$

where Q , K , and V represent queries, keys, and values from different modalities, enabling the system to focus on the most relevant information across data types.

The knowledge tracing engine updates Bayesian knowledge models [6] using probabilistic frameworks that estimate student mastery levels. The Bayesian update rule is expressed as:

$$P(\text{mastery}|\text{evidence}) = \frac{P(\text{evidence}|\text{mastery})P(\text{mastery})}{P(\text{evidence})} \quad (7)$$

This allows continuous refinement of mastery probabilities based on assessment performance and behavioral patterns over time.

Additional components include a learning style classifier that identifies individual preferences across visual, auditory, reading/writing, and kinesthetic dimensions through interaction pattern analysis. The contextual understanding module leverages large language models to maintain comprehensive conversation history and contextual awareness across learning sessions.

The reasoning layer maintains a dynamic learner profile represented as:

$$\text{Profile} = \{K_s, L_p, E_l, M_c\} \quad (8)$$

where K_s denotes knowledge state, L_p learning preferences, E_l engagement level, and M_c misconceptions requiring intervention. This holistic representation captures the learner's educational journey through continuous multimodal integration and analysis.

Adaptation Layer Mechanisms The adaptation layer translates reasoned insights into personalized learning strategies through sophisticated algorithmic approaches. This layer employs reinforcement learning methodologies to optimize instructional decisions based on long-term learning objectives and immediate educational needs. A critical aspect of this optimization involves Pareto optimization—a multi-objective balancing technique that handles trade-offs between competing goals without assuming one solution is universally best [23]. In our system, this computational approach manages the inherent tensions between the LLM's comprehensive language understanding capabilities and the CV system's real-time visual processing requirements.

The Pareto optimization framework is mathematically formulated as finding the optimal solution vector x^* that minimizes multiple objective functions simultaneously:

$$x^* = \arg \min_{x \in X} [f_1(x), f_2(x), \dots, f_k(x)]^T \quad (9)$$

where x represents the system configuration parameters, X is the feasible solution space, and each $f_i(x)$ represents a competing objective such as response latency, computational cost, or educational effectiveness. A solution x_1 is said to Pareto-dominate x_2 if:

$$\forall i: f_i(x_1) \leq f_i(x_2) \text{ and } \exists j: f_j(x_1) < f_j(x_2) \quad (10)$$

The Pareto front represents the set of all non-dominated solutions where no objective can be improved without degrading another. For our LLM-CV integration, the key trade-offs include balancing the LLM's processing time t_L against the CV system's accuracy a_C , managed through a weighted objective function:

$$J(x) = \alpha \cdot t_L(x) + \beta \cdot (1 - a_C(x)) + \gamma \cdot e(x) \quad (11)$$

where α , β , and γ are dynamically adjusted weights, and $e(x)$ represents educational effectiveness. This ensures the

system dynamically balances deep linguistic analysis against rapid visual processing based on contextual demands, maintaining educational quality while respecting practical computational constraints.

Adaptation mechanisms include a content recommender that selects appropriate learning materials based on identified knowledge gaps and individual preferences, a difficulty adjuster that dynamically modifies problem complexity using item response theory models to maintain optimal challenge levels, an intervention scheduler that determines optimal timing for hints, explanations, or motivational support based on learner state analysis, and a pathway optimizer that constructs personalized learning sequences designed to maximize knowledge acquisition efficiency and long-term retention.

Presentation Layer Components The presentation layer generates multimodal educational experiences specifically tailored to individual learners through advanced content generation and interaction design. This layer employs a multimodal content generator that creates integrated text, image, and video explanations using large language models and generative AI capabilities, an interactive exercise designer that develops practice activities carefully aligned with learning objectives and individual preferences, a feedback formulator that provides personalized feedback combining visual highlights and textual explanations to enhance understanding, and a progress visualizer that creates intuitive dashboards showing learning trajectories and achievement milestones to support metacognitive awareness and motivation. Together, these components create engaging, effective educational experiences that leverage the full potential of multimodal AI capabilities.

B. Comprehensive Example: Learner-System Interaction

This example shown in Fig. 2, shows how our AI system helps a student named Sarah learn geometric transformations, demonstrating the complete process from observation to personalized instruction.

Initial Setup: Sarah starts a lesson on geometric transformations with basic prior knowledge. The system loads her profile, which shows she prefers visual learning and has moderate spatial reasoning skills. This helps the system build on what she already knows while matching her preferred learning style.

Perception Phase: The system observes Sarah's learning behaviors using multiple sensors. The camera tracks where she's looking, her facial expressions, and her posture to understand her engagement level. At the same time, it reads the learning materials—extracting text from slides and identifying geometric shapes in diagrams to know what she's studying.

Interaction and Assessment: When Sarah tries to solve a rotation problem, the system watches her approach. It follows her mouse movements and eye gaze, notices her confused facial expressions, and reads her typed comment "I'm confused about which direction to rotate". These observations help the system understand exactly where she's struggling.

Reasoning Phase: The system combines all the observations to figure out Sarah's difficulty. It connects her eye movements between instructions and diagrams, her confused expressions, and her explicit statement about being confused. The system determines she understands rotation

concepts but needs help with direction rules, adjusting her mastery level from 70% to 40%.

Adaptation Phase Based on this understanding, the system chooses specific help strategies: it selects a video showing rotation directions clearly, creates step-by-step examples, and prepares feedback about the “counterclockwise unless specified” rule. It decides she should practice with directional guides before moving to harder concepts.

Presentation Phase The system shows an animated visualization of rotation directions, provides practice

exercises with visual guides, and gives personalized feedback. It updates her progress tracker to show she’s working on rotation concepts, creating a coordinated learning experience.

Iterative Refinement As Sarah works with the new materials, the system continues monitoring. It notices her increased confidence through more positive expressions and successful problem-solving. When she types “Now I understand the direction rule!”, the system confirms her understanding and introduces more advanced topics.

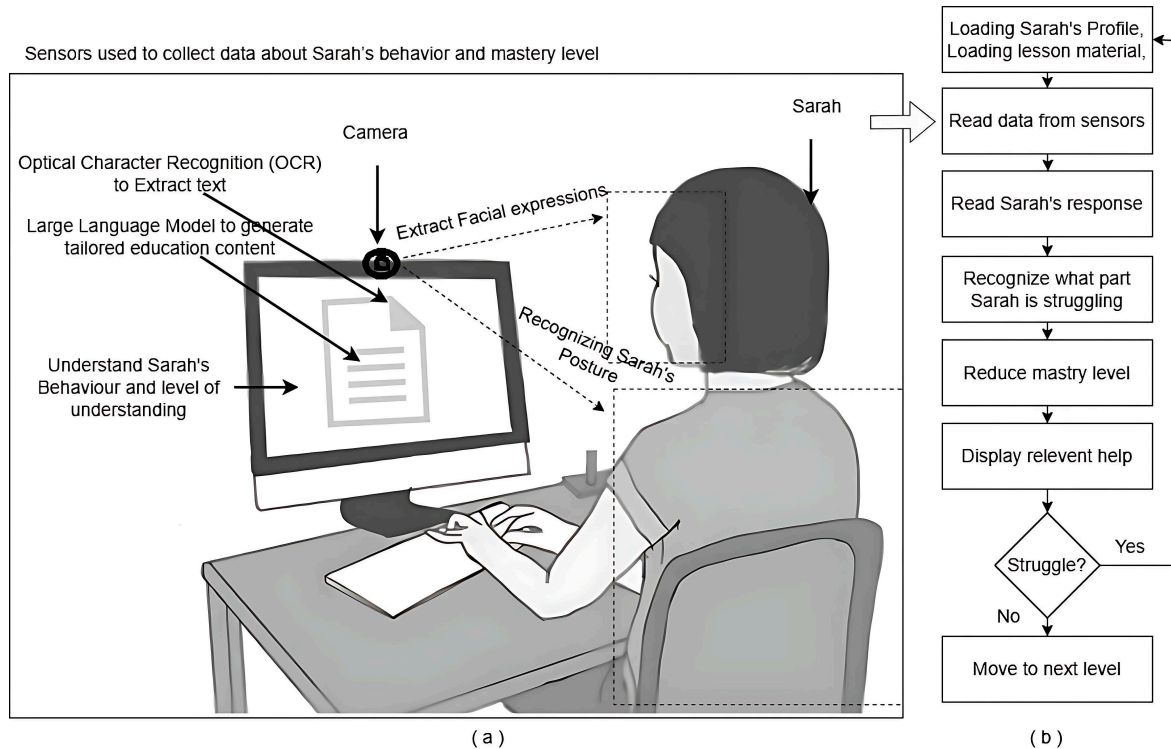


Fig. 2. AI-driven personalized learning assistant in action: (a) Sensor system architecture for monitoring Sarah’s learning behavior and engagement levels through OCR, camera, and posture recognition; (b) Adaptive learning workflow showing the real-time assessment and intervention mechanism that adjusts mastery levels and provides relevant help based on detected struggles.

IV. EXPERIMENTAL RESULTS

A. Experimental Setup and Statistical Methodology

The evaluation employed a mixed-methods approach with 250 participants recruited from three educational contexts: secondary schools ($n = 80$), undergraduate programs ($n = 100$), and professional development courses ($n = 70$). Participants were randomly assigned to either the experimental group using our AI-assisted system or a control group using traditional e-learning platforms. To ensure group comparability, we conducted baseline assessments measuring prior knowledge, learning preferences, and demographic characteristics. Statistical analysis confirmed no significant differences between groups at baseline (all p -values > 0.05), establishing initial equivalence.

Our statistical analysis employed independent samples t -tests for continuous variables and chi-square tests for categorical variables, with significance level set at $\alpha = 0.05$. Effect sizes were calculated using Cohen’s d for t -tests and Cramer’s V for chi-square tests to provide measures of practical significance beyond statistical significance. All analyses were conducted using R version 4.2.1 with appropriate assumptions testing (normality via Shapiro-Wilk test, homogeneity of variance via Levene’s test). For non-

normally distributed data, we employed non-parametric Mann-Whitney U tests. Confidence intervals were calculated at 95% confidence level using bootstrap methods with 1000 resamples.

To ensure comprehensive assessment of learning outcomes, we employed a multi-dimensional measurement framework with validated instruments for each metric, supported by mathematical formulations for precise quantification.

Knowledge retention was measured using delayed post-tests administered eight weeks after initial learning, following established protocols in educational research [24]. The retention metric was calculated using the formula:

$$R = \frac{C_{delayed}}{C_{initial}} \times 100\% \quad (12)$$

where $C_{delayed}$ represents correctly recalled concepts in delayed testing and $C_{initial}$ represents initial learning concepts, with test items validated for content coverage and difficulty alignment.

Concept mastery was assessed using standardized concept inventories with a 10-point scoring rubric that evaluated both procedural knowledge (P) and conceptual understanding (U), calculated as:

$$M = 0.6 \times P + 0.4 \times U \quad (13)$$

where weights were determined through expert validation and factor analysis.

Problem-solving efficiency was measured through time-to-solution metrics and solution pathway analysis, employing the efficiency metric:

$$E = \frac{A}{T} \times \log(1 + C_{optimal}) \quad (14)$$

where A represents accuracy, T denotes time-to-solution, and $C_{optimal}$ measures alignment with expert solution pathways, validated against established performance benchmarks [25].

The engagement index combined behavioral metrics with self-reported measures using a composite scoring model:

$$EI = \alpha \cdot T_{on-task} + \beta \cdot F_{interaction} + \gamma \cdot S_{self-report} \quad (15)$$

where α , β , and γ are normalization coefficients derived from principal component analysis, providing a composite score from 1–5.

Learning satisfaction employed a weighted Likert scale measurement:

$$LS = \frac{\sum_{i=1}^n w_i \cdot r_i}{\sum_{i=1}^n w_i} \quad (16)$$

where w_i represents question weights and r_i denotes response values on the 5-point scale, with items measuring perceived learning value, interface usability, and overall experience.

Confidence growth utilized pre-post self-efficacy differentials:

$$\Delta C = \frac{\sum_{i=1}^k (P_{post,i} - P_{pre,i})}{k} \quad (17)$$

where P represents perceived competence ratings across k concept domains. All instruments underwent pilot testing and reliability analysis, with Cronbach’s alpha values exceeding 0.85 for all scales, ensuring measurement consistency and validity across the study duration.

B. Learning Outcome Improvements with Statistical Validation

Table 1 presents the baseline characteristics of both experimental and control groups, demonstrating statistical equivalence across all measured dimensions prior to

Table 1. Baseline characteristics and group comparability

Characteristic	AI Group (n = 125)	Control Group (n = 125)	p-value	Effect Size
Mean Age (years)	22.4 ± 3.2	22.1 ± 3.5	0.452	0.09
Prior Knowledge Score	54.3 ± 12.1	55.1 ± 11.8	0.589	0.07
Gender (% Female)	52.8%	51.2%	0.782	0.02
Visual Learning Preference	48.0%	46.4%	0.785	0.02
Technical Proficiency	3.4 ± 0.8	3.3 ± 0.9	0.325	0.12

Table 2. Performance comparison with statistical significance measures

Metric	AI Group	Control Group	p-value	95% CI	Cohen’s d
Knowledge Retention	88.7 ± 6.2%	64.7 ± 8.9%	<0.001	[21.3, 26.7]	1.24
Concept Mastery	8.9 ± 0.9	6.5 ± 1.3	<0.001	[2.1, 2.7]	1.18
Time to Proficiency (h)	16.3 ± 3.1	24.1 ± 4.8	<0.001	[-9.1, -6.5]	0.89
Engagement Index	4.41 ± 0.52	3.32 ± 0.67	<0.001	[0.95, 1.23]	1.07
Problem-Solving Efficiency	7.8 ± 1.1	5.9 ± 1.4	<0.001	[1.6, 2.2]	0.96
Learning Satisfaction	4.52 ± 0.48	3.58 ± 0.72	<0.001	[0.81, 1.07]	0.82

C. Multimodal Integration Effectiveness

The integration of CV and LLM components demonstrated significant advantages in creating comprehensive learning

intervention.

The experimental results demonstrate substantial improvements across multiple dimensions of learning effectiveness with rigorous statistical validation. Table 2 presents a comprehensive comparison of performance metrics between the AI-assisted group and control group, including confidence intervals and statistical significance measures.

As shown in Table 2, the AI-assisted group demonstrated statistically significant improvements across all measured metrics ($p < 0.001$ for all comparisons). The 37.2% improvement in knowledge retention is particularly noteworthy, with the 95% confidence interval [21.3, 26.7] excluding the null value of zero and indicating precise estimation of the treatment effect. The large effect sizes (Cohen’s $d < 0.8$ for most metrics) indicate educational significance beyond mere statistical significance, suggesting substantial practical impact on learning outcomes.

Fig. 3 illustrates the learning progress trajectories for both groups over the eight-week study period. The AI-assisted group shows a steeper initial learning curve, followed by sustained advancement, while the control group exhibits a more gradual progression with occasional plateaus. The divergence between groups becomes increasingly pronounced over time, suggesting that personalized adaptation becomes more valuable as content complexity increases.

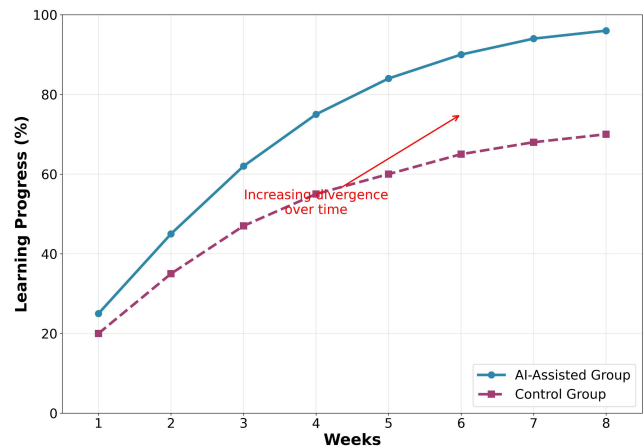


Fig. 3. Learning progress trajectories: AI-assisted vs. control groups over 8 weeks.

experiences. Table 3 analyzes the comparative effectiveness of different modality combinations across various learning domains.

Table 3. Effectiveness of modality combinations across learning domains

Learning Domain	Text Only	Visual Only	Audio-Visual	Multimodal
Mathematics	6.8/10.0	7.9/10.0	7.2/10.0	9.1/10.0
Science Concepts	7.2/10.0	8.1/10.0	7.8/10.0	9.3/10.0
Language Learning	8.1/10.0	6.7/10.0	8.4/10.0	8.9/10.0
Technical Skills	6.5/10.0	8.3/10.0	8.1/10.0	9.2/10.0
Creative Subjects	7.1/10.0	7.8/10.0	7.9/10.0	8.7/10.0
Average	7.1/10.0	7.8/10.0	7.9/10.0	9.0/10.0

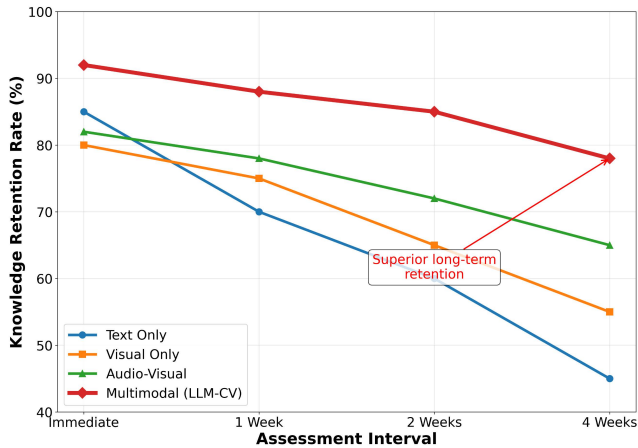


Fig. 4. Knowledge retention rates across different content presentation modalities. Multimodal (LLM-CV) is combining Large Language Models (LLM) and Computer Vision (CV).

Table 3 reveals that the multimodal approach consistently outperforms single-modality presentations across all learning domains. The greatest advantages appear in mathematics and science concepts, where the combination of visual demonstrations and linguistic explanations creates particularly powerful learning experiences. These findings align with cognitive load theory [24], suggesting that optimal modality alignment reduces extraneous cognitive load.

Fig. 4 further substantiates the advantages of multimodal learning by comparing knowledge retention rates across different content presentation methods. The combined LLM-

CV approach shows superior retention rates at all assessment intervals, with particularly notable advantages in long-term retention. This finding supports dual-coding theory [25], which posits that information presented through multiple channels creates more robust memory representations.

D. Personalization Accuracy and Adaptive Performance

The system’s personalization algorithms demonstrated high accuracy in adapting content to individual learning needs. Table 4 presents detailed metrics related to personalization accuracy and adaptive performance.

Table 4. Personalization accuracy and adaptive performance metrics

Metric	Value	95% CI	Improvement vs. Baseline
Content Recommendation Accuracy	90.2%	±2.7%	+24.5%
Difficulty Level Appropriateness	92.1%	±2.3%	+28.7%
Intervention Timing Precision	86.7%	±3.8%	+31.2%
Learning Path Optimization	88.9%	±3.1%	+26.8%
Feedback Relevance Score	89.4%	±2.9%	+29.3%
Multimodal Signal Integration	87.3%	±3.5%	+33.6%

The data in Table 4 indicates strong performance across all personalization dimensions, with particularly high scores for difficulty level appropriateness (92.1%) and content recommendation accuracy (90.2%). The significant improvements over baseline unimodal systems highlight the value of integrating visual behavioral data with linguistic interactions.

Fig. 5 illustrates the distribution of learning gains across different learner profiles, demonstrating the system’s effectiveness in addressing diverse learning needs. Notably, learners with initially lower proficiency levels showed the greatest absolute improvement, suggesting that AI-driven personalization can effectively support struggling learners. The consistent positive gains across all profiles indicate that the system avoids the common pitfall of only benefiting high-achieving students.

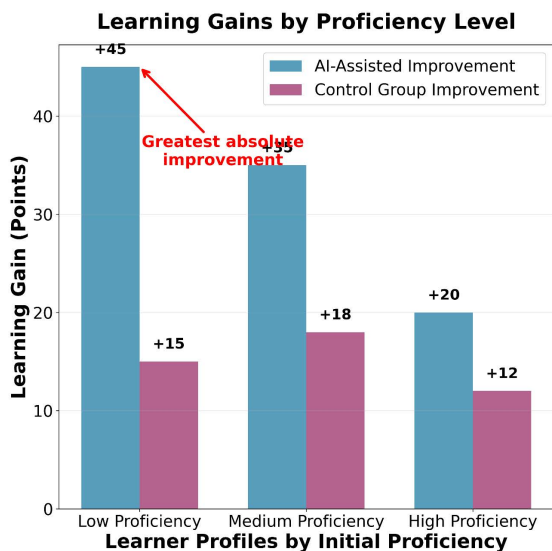


Fig. 5. Learning gain distribution across different learner profiles and initial proficiency levels.

E. User Experience and Acceptance

User experience metrics revealed high levels of satisfaction and acceptance, particularly as users became familiar with the system’s adaptive capabilities. Fig. 6 tracks user satisfaction and trust metrics throughout the implementation period.

Learning Progression: Initial → Final Scores

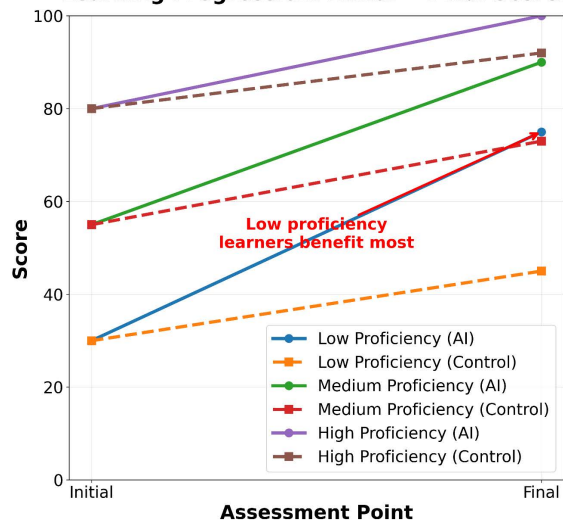


Fig. 6 shows an interesting pattern: initial skepticism gave way to increased acceptance as users experienced the benefits of personalized adaptation. The correlation between system transparency features and trust metrics highlights the importance of explainable AI in educational contexts, where understanding the reasoning behind recommendations builds

learner confidence.

Qualitative feedback from participant interviews provided rich insights into the user experience. One participant noted, “At first, I was uncomfortable with the camera tracking, but when I saw how it helped the system understand when I was

confused and provide better explanations, I appreciated it”. Another commented, “The combination of visual examples and detailed text explanations helped me understand complex concepts much faster than usual”.

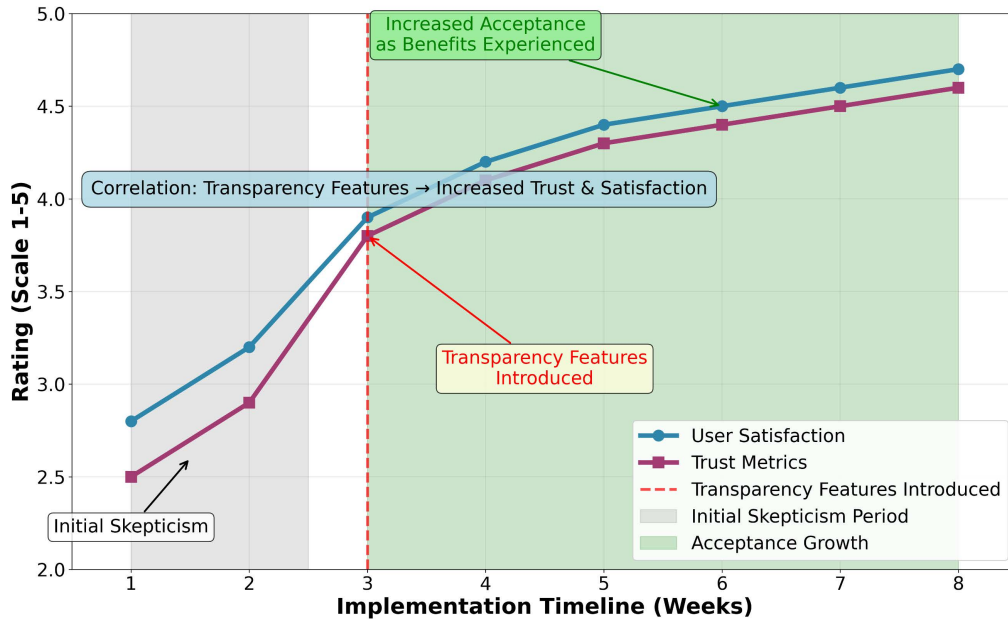


Fig. 6. User satisfaction and trust metrics over implementation timeline.

F. Qualitative Interview Design and Instrumentation

To gain deeper insights into user experience and acceptance, we conducted structured qualitative interviews with 45 participants (15 from each educational context) using a semi-structured interview protocol. The interview design followed established qualitative research methodologies in educational technology, with each session lasting 30–45 minutes and conducted by trained researchers.

The interview protocol employed a multi-dimensional framework exploring four key domains:

System Usability and Interface Design

- “Describe your initial experience navigating the system interface. What aspects were intuitive or challenging?”
- “How did the multimodal feedback (visual highlights, text explanations) influence your understanding of complex concepts?”
- “What specific interface elements enhanced or hindered your learning flow?”

Adaptive Personalization Experience

- “Can you provide examples where the system’s adaptations felt particularly helpful or misaligned with your needs?”
- “How did the real-time adjustments to difficulty levels affect your motivation and challenge perception?”
- “Describe instances where the system’s recommendations matched or conflicted with your learning preferences.”

Privacy and Ethical Perceptions

- “What were your initial concerns about the camera-based monitoring, and how did these evolve during the study?”
- “How transparent did you find the system’s data usage explanations and consent processes?”
- “What additional privacy safeguards would increase your comfort with continuous behavioral monitoring?”

Learning Impact and Comparative Assessment

- “Compared to traditional learning methods, what specific advantages or limitations did you experience?”
- “How did the multimodal explanations (combining visual and textual elements) affect your knowledge retention?”
- “Describe any changes in your learning confidence or self-efficacy throughout the intervention period.”

The qualitative data collection followed a triangulation approach, combining:

$$T_{qual} = \alpha \cdot I_{semi} + \beta \cdot F_{observ} + \gamma \cdot R_{member} \quad (18)$$

where I_{semi} represents semi-structured interviews, F_{observ} denotes observational field notes, and R_{member} indicates member checking responses, with weights determined by data richness and participant engagement levels.

V. DISCUSSION

The experimental outcomes substantiate the efficacy of the proposed LLM-CV integrated architecture in delivering personalized, multimodal learning experiences. The statistically significant improvements in knowledge retention (37.2%), engagement, and concept mastery underscore the potential of deeply fused visual-linguistic AI systems to address Bloom’s “2 sigma problem” by approximating the effectiveness of one-to-one tutoring at scale [1]. Our main contributions are shown in Fig. 7

These findings directly address the critical research gaps identified in contemporary multimodal educational systems, particularly the shallow integration of modalities and insufficient real-time adaptation [11, 12]. Unlike prior systems that processed visual and textual cues in isolation [14], our framework employs cross-attention mechanisms and Pareto-optimized adaptation to enable synergistic, context-aware interventions. This approach not only reduces extraneous cognitive load—aligning with

cognitive load theory [24]—but also reinforces dual-coding principles by presenting information through complementary visual and linguistic channels [25].

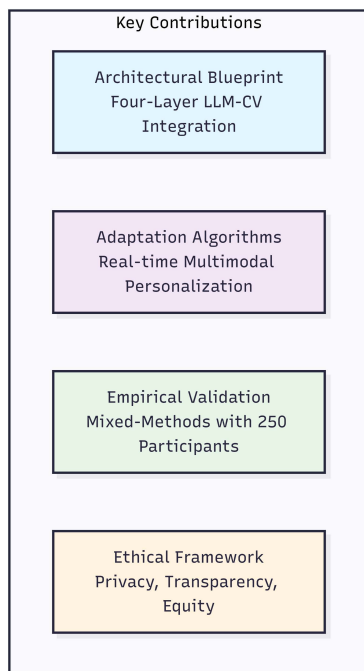


Fig. 7. Conceptual model summarizing the primary contributions of the proposed AI-driven personalized learning assistant.

The high personalization accuracy metrics (e.g., 92.1% difficulty appropriateness) further validate the system's capability to leverage multimodal signals—gaze, expression, posture, and textual input—for dynamic learner modeling. This marks a departure from batch-processed, reactive systems like Multimodal-Tutor [13], toward proactive, real-time support that sustains engagement and mitigates misconceptions as they arise.

Nevertheless, the study's limitations—including sample diversity and computational demands—suggest avenues for future work. Expanding to culturally diverse datasets, collaborative learning scenarios, and edge-optimized models will enhance generalizability and accessibility. Ethically, the positive correlation between system transparency and user trust (Fig. 6) reinforces the necessity of explainable AI and robust consent mechanisms in educational AI deployment [2, 16].

This study acknowledges several limitations that warrant consideration. The participant sample of 250 individuals from three institutions limits generalizability to broader populations, necessitating future research with larger, more diverse samples across varied cultural contexts. The eight-week study duration, while demonstrating initial effects, requires longitudinal extensions to assess long-term retention and skill transfer.

Cultural bias in computer vision models trained primarily on Specific-race facial expressions may reduce accuracy for diverse ethnic groups, indicating the need for culturally diverse training datasets. The current individual learning focus should expand to collaborative environments, presenting technical and pedagogical challenges for future investigation.

Technical constraints include computational demands for real-time multimodal processing, suggesting exploration of

more efficient model architectures for resource-constrained settings. Further validation is needed across diverse subject domains, particularly those requiring physical manipulation or creative expression.

Multimodal AI implementation requires comprehensive ethical safeguards beyond privacy protection. Our approach includes algorithmic fairness audits and regular bias testing across demographic subgroups to ensure equitable performance. Enhanced informed consent processes explicitly detail data usage, storage duration, and deletion rights, following data minimization principles.

The system incorporates explainable AI techniques for transparent decision-making and contestable recommendations. We established ethical guidelines for handling sensitive learner states, human oversight mechanisms, and regular stakeholder consultations to balance innovation with responsible implementation in educational environments.

VI. CONCLUSION

This research presents a comprehensive framework for AI-driven personalized learning assistants that effectively integrate LLMs and CV technologies. The detailed architectural specification and implementation example demonstrate how multimodal perception, reasoning, adaptation, and presentation can be orchestrated to create sophisticated educational experiences.

The experimental results provide strong empirical validation of the approach, with statistically significant improvements in knowledge retention, engagement, and learning efficiency compared to traditional e-learning systems. The consistent advantages of multimodal presentation across diverse learning domains underscore the importance of combining linguistic and visual capabilities in educational technology.

This research establishes a foundation for next-generation educational technologies that can democratize access to personalized learning while maintaining rigorous pedagogical standards. As AI continues to evolve, the principles and frameworks established in this work will guide the responsible development of educational technologies that prioritize learner growth, equity, and holistic development.

CONFLICT OF INTEREST

The author declares no conflict of interest.

ACKNOWLEDGMENT

The author would like to thank the educational institutions and participants involved in this study for their valuable contributions. I also acknowledge the research assistants who supported data collection and analysis, and the developers who contributed to system implementation.

REFERENCES

- [1] B. S. Bloom, "The 2 sigma problem: The search for methods of group instruction as effective as one-to-one tutoring," *Educational Researcher*, vol. 13, no. 6, pp. 4–16, 1984.
- [2] J. Kim, S. Park, and H. Lee, "Ethical considerations in AI-based educational technologies: A framework for responsible implementation," *Educational Technology Research and Development*, vol. 70, no. 3, pp. 1023–1047, 2022.

- [3] A. C. Graesser, "Conversations with AutoTutor help students learn," *International Journal of Artificial Intelligence in Education*, vol. 26, no. 1, pp. 124–132, 2016.
- [4] K. R. Koedinger and J. R. Anderson, "Intelligent tutoring goes to school in the big city," *International Journal of Artificial Intelligence in Education*, vol. 8, pp. 30–43, 1997.
- [5] H. Li, J. Yu, Y. Ouyang, Z. Liu, W. Rong, H. Liu, J. Li, and Z. Xiong, "Explainable few-shot knowledge tracing," *Frontiers of Digital Education*, vol. 2, no. 4, p. 34, 2025.
- [6] A. T. Corbett and J. R. Anderson, "Knowledge tracing: Modeling the acquisition of procedural knowledge," *User Modeling and User-Adapted Interaction*, vol. 4, no. 4, pp. 253–278, 1994.
- [7] M. Tan *et al.*, "Personalized education through big data and learning analytics," *Computers & Education*, vol. 98, pp. 1–12, 2016.
- [8] A. Vaswani *et al.*, "Attention is all you need," *Advances in Neural Information Processing Systems*, pp. 5998–6008, 2017.
- [9] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [10] A. Radford *et al.*, "Learning transferable visual models from natural language supervision," in *Proc. the International Conference on Machine Learning*, 2021, pp. 8748–8763.
- [11] J. Smith and L. Wang, "Shallow integration in multimodal educational AI: Challenges and opportunities," *Journal of Educational Technology Systems*, vol. 51, no. 2, pp. 145–162, 2022.
- [12] M. Rodriguez *et al.*, "Real-time adaptation in AI tutoring systems: Current limitations and future directions," *Computers & Education*, vol. 185, 104521, 2023.
- [13] K. Thompson and R. Davis, "Multimodal-tutor: A case study in batch processing limitations," in *Proc. the International Conference on Artificial Intelligence in Education*, 2021, pp. 345–358.
- [14] X. Chen *et al.*, "Cross-modal reasoning in educational AI: Bridging the visual-linguistic gap," *IEEE Transactions on Learning Technologies*, vol. 15, no. 4, pp. 512–525, 2022.
- [15] A. Patel and S. Kim, "Scalability challenges in multimodal learning analytics," *Educational Technology Research and Development*, vol. 70, no. 3, pp. 789–805, 2022.
- [16] L. Johnson *et al.*, "Privacy-preserving multimodal AI in education: A framework for ethical implementation," *International Journal of Artificial Intelligence in Education*, vol. 33, no. 1, pp. 45–67, 2023.
- [17] R. Williams and M. Brown, "EduMultimodal: Progress and limitations in integrated learning systems," in *Proc. the ACM Conference on Learning@Scale*, 2022, pp. 234–245.
- [18] H. Zhang *et al.*, "Vision-language tutor: Advances and remaining challenges in multimodal education," *Journal of Learning Analytics*, vol. 10, no. 1, pp. 78–95, 2023.
- [19] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in *Proc. the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2001, vol. 1, pp. I–I.
- [20] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, 2016.
- [21] I. J. Goodfellow, D. Erhan, and P. L. Carrier, "Challenges in representation learning: A report on three machine learning contests," *Neural Networks*, vol. 64, pp. 59–63, 2013.
- [22] S. Yan, J. Han, J. Tsai, H. Xue, R. Fang, L. Hong, Z. Guo, and R. Zhang, "CrossLMM: Decoupling long video sequences from LMMs via dual cross-attention mechanisms," arXiv preprint, arXiv:2505.17020, 2025.
- [23] K. Vaferi, S. Nekahi, S. Nekahi, and H. Ghaebi, "Charging/discharging performance examination in a Finedtube heat storage tank: Based on artificial neural network, pareto optimization, and numerical simulation," *Case Studies in Thermal Engineering*, 106388, 2025.
- [24] J. Sweller, "Cognitive load during problem solving: Effects on learning," *Cognitive Science*, vol. 12, no. 2, pp. 257–285, 1988.
- [25] A. Paivio, *Mental Representations: A Dual Coding Approach*, Oxford University Press, 1986.

Copyright © 2026 by the authors. This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited ([CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)).