

# Incorporating brain-inspired mechanisms for multimodal learning in artificial intelligence

Xiang He<sup>1,\*</sup>, Dongcheng Zhao<sup>1,2,\*</sup>, Yang Li<sup>1</sup>, Qingqun Kong<sup>1†</sup>, Xin Yang<sup>3†</sup>, Yi Zeng<sup>1,2,4†</sup>

<sup>1</sup>Brain-inspired Cognitive AI Lab, Institute of Automation, Chinese Academy of Sciences, Beijing, China

<sup>2</sup>Center for Long-term AI, Beijing, China

<sup>3</sup>CAS Key Laboratory of Molecular Imaging, Institute of Automation, Chinese Academy of Sciences, Beijing, China

<sup>4</sup>Key Laboratory of Brain Cognition and Brain-inspired Intelligence Technology, Chinese Academy of Sciences, Shanghai, China

{hexiang2021, zhaodongcheng2016, liyang2019, qingqun.kong, xin.yang, yi.zeng}@ia.ac.cn

## Abstract

*Multimodal learning significantly enhances the perceptual capabilities of cognitive intelligent systems by integrating information from different sensory modalities. However, existing multimodal fusion researches in the field of artificial intelligence typically assume static integration of modal information, not yet fully incorporating the key dynamic mechanisms of multimodal integration found in the brain. Specifically, when processing multisensory information, the brain exhibits an inverse effectiveness phenomenon, wherein weaker unimodal cues yield stronger multisensory integration benefits; conversely, when individual modal cues are stronger, the effect of modal fusion is relatively diminished. This mechanism enables biological systems to achieve robust cognition even in environments with scarce or noisy perceptual cues. Inspired by this biological inverse effectiveness mechanism in multimodal integration, we explore the intrinsic relationship between multimodal output and information from individual modalities, proposing an inverse effectiveness driven multimodal fusion (IEMF) strategy. By incorporating this inverse effectiveness-driven multimodal fusion strategy into neural network architectures, we achieve not only more efficient multimodal integration with significantly improved model performance, but also substantial computational efficiency gains—demonstrating up to 50% reduction in computational cost across diverse fusion methods. We conduct extensive experiments on audio-visual classification, audio-visual continual learning, and audio-visual question answering tasks to validate the effectiveness of our proposed method. The experimental results consistently demonstrate that our proposed method performs excellently in these multimodal tasks. Furthermore, to verify the universality and*

*generalization capability of our method, we also conduct experiments on two widely used network models in artificial intelligence—Artificial Neural Networks (ANN) and Spiking Neural Networks (SNN)—with results showing good adaptability of the method to both network types. Our research emphasizes the potential of incorporating biologically inspired neural mechanisms into multimodal neural networks and provides promising new directions and perspectives for the future research and development of multimodal artificial intelligence. The code is publicly available at <https://github.com/Brain-Cog-Lab/IEMF>.*

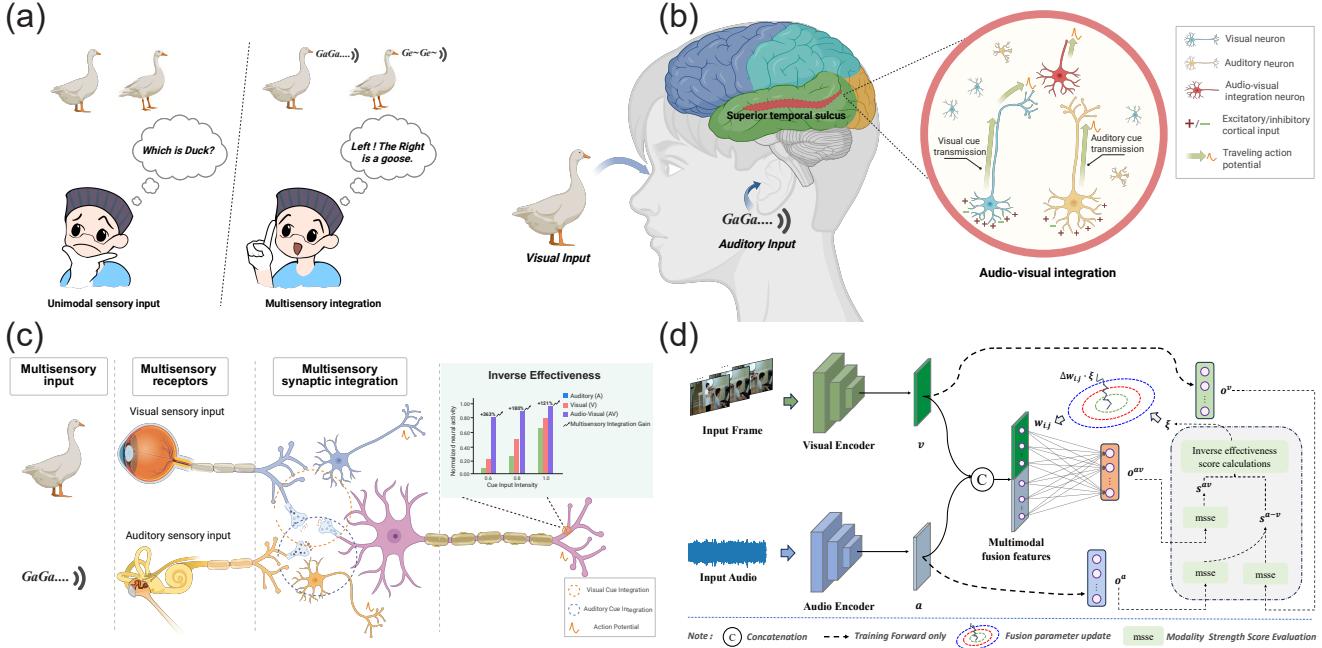
## 1. Introduction

In natural environments, we typically need to process cues from multiple senses simultaneously to comprehensively construct an understanding of the same concept. For example, understanding the concept of "beach" involves not only visual information (yellow sand, blue sea) but also relies on auditory (sound of waves) and tactile (texture of sand) sensory information. Compared to unimodal information, multimodal information provides richer and more comprehensive representational capacity [13, 35]. As shown in Fig. 1(a), multimodal integration not only enhances information expressiveness but also effectively reduces uncertainty in single-modal information. This mechanism of multimodal information integration is not only the foundation of biological perception but has also become one of the core challenges in multimodal learning in artificial intelligence. As information environments become increasingly complex, traditional unimodal learning methods struggle to handle complex and dynamic real-world scenarios. Consequently, neural networks have incorporated multimodal information processing strategies to achieve more robust and efficient information representation. Mul-

---

\*Equal Contribution

†Corresponding Author



**Figure 1. Illustration of multisensory integration and the role of inverse effectiveness in IEMF (Inverse Effectiveness driven Multimodal Fusion).** (a) Comparison between unimodal sensory input and multisensory integration: integrating visual and auditory cues reduces ambiguity and uncertainty compared to relying on a single modality. (b) Neural basis of audiovisual integration in the human brain, focusing on the superior temporal sulcus (STS) where visual and auditory inputs converge onto multisensory neurons. (c) Biological principle of inverse effectiveness: multisensory integration is strengthened when unimodal signals are weak. Visual and auditory stimuli are processed through distinct sensory pathways and converge at multisensory synapses. The inset illustrates the inverse relationship between unimodal strength and integrative gain. (d) The proposed inverse effectiveness driven multimodal fusion strategy inspired by biological multisensory fusion mechanisms. Visual and auditory inputs are processed by respective encoders, fused via a dynamic fusion module regulated by inverse effectiveness principles, and evaluated using modality strength score estimation. The fusion module weights are dynamically adjusted according to the computed scores. (This figure was created with <https://BioRender.com>.)

timodal neural networks are widely applied in tasks such as multimodal fusion [21, 32, 37, 41, 53], multimodal emotion recognition [10, 28], and audio-visual speech recognition [24, 31, 37, 52]. Nevertheless, brain-inspired algorithms remain in the developmental stage, and many biological characteristics and mechanisms have not yet been fully utilized, offering enormous potential and new challenges for the further development of neural network models.

Neurobiological research indicates that vision and audition are the two primary pathways through which humans acquire external information, and their integration significantly enhances perceptual benefits [6, 12]. This paper therefore focuses primarily on the integration of visual and auditory modal inputs. For visual and auditory information from a common source, the brain has specialized regions responsible for both unimodal information processing and multisensory integration [29]. After visual and auditory information are received through their respective receptors, features are hierarchically extracted through visual and auditory pathways before being transmitted to multisensory integration brain regions. Relevant studies show that au-

diovisual information integration occurring in the cerebral cortex is closely associated with regions such as the superior temporal sulcus [8, 30, 34, 42, 47, 49], posterior parietal cortex [3, 39], and prefrontal cortex [4, 7]. Figure 1(b) illustrates the convergence and integration process of co-sourced audiovisual signals in the superior temporal sulcus as an example. In this multisensory integration brain region, visual and auditory cues are transmitted through different neural pathways, ultimately converging onto common multisensory integration neurons, facilitating cross-modal information integration and perceptual decision-making.

Brain audiovisual information integration exhibits many interesting mechanisms, with inverse effectiveness being particularly noteworthy. [8] found that by separately presenting audiovisual combined speech signals and their unimodal information, and conducting cross-modal comparisons, the left superior temporal sulcus (STS) demonstrated the most significant cross-modal integration benefits under conditions where unimodal signals were weakest. The inverse effectiveness mechanism indicates that during multisensory information integration, when unimodal cues

are weaker, the effect of multisensory integration is relatively stronger; conversely, when individual modal cues are stronger, the effect of modal fusion is relatively diminished, though multisensory integration responses still exceed the activation response of either single modality [14, 45]. Figure 1(c) illustrates this phenomenon, where inverse effectiveness reflects higher sensitivity to weaker modalities in multimodal integration brain regions, typically manifested as enhanced information integration. This mechanism enables biological systems to enhance perceptual accuracy and stability by strengthening multimodal integration when the quality of information from a single modality is poor.

Inspired by the biological principle of inverse effectiveness, our work reconsiders how multimodal fusion should adapt to variations in unimodal input quality, particularly under dynamic and complex environments. Most existing multimodal fusion methods focus on maximizing information interaction between modalities, yet often overlook the dynamic relationship between the fused output and the respective contributions of each modality. This limitation stems from existing methods typically presetting modal interactions as static fixed patterns, failing to fully consider that different modalities' information contributions should flexibly adjust as environmental conditions change. Take audiovisual perception as an example: when environmental noise significantly degrades the quality of auditory input, traditional fusion strategies typically retain fixed fusion weights and cannot adaptively modify the cooperative interaction between modalities in accordance with signal degradation, thereby constraining the overall perceptual performance of the system. The phenomenon of inverse effectiveness inspires the insight that an efficient multimodal integration mechanism should actively enhance the responsiveness of the fusion module when the quality of a single modality degrades, so that the system can obtain more compensatory information from other modalities. Based on this insight, we propose that fusion strength should dynamically respond to modality-specific quality fluctuations, that is, the learning rate of the fusion module should be adaptively modulated according to the reliability of unimodal signals, thereby enabling more robust and flexible multimodal perception in complex and evolving environments.

Based on the aforementioned neural mechanisms of multimodal fusion and biological inspiration, this paper adopts deep neural networks as the foundational framework for multimodal perceptual learning, focusing on exploring the cooperative integration process of visual and auditory information. We propose an inverse effectiveness driven multimodal fusion (IEMF) strategy to enable a more fine-grained fusion mechanism. By quantifying the relationship between the strength of unimodal inputs and the signal strength of multimodal fusion outputs, we adaptively modulate the update rate of the fusion module's weights. Specifically, we

introduce an inverse effectiveness coefficient into the back-propagation process, such that the fusion module accelerates its parameter updates in response to weak unimodal signals to enhance fusion strength, while suppressing updates when unimodal signals are strong, thereby reducing over-reliance on fusion. This design realizes a biologically inspired principle of "weak modality, strong fusion," ensuring that the integration process neither over-depends on a single sensory pathway nor overlooks potentially informative sources. Consequently, the method effectively improves overall perceptual accuracy and model robustness. Beyond improving overall perceptual accuracy and model robustness, our approach also demonstrates significant computational efficiency gains—reducing training costs by up to 50% while maintaining superior performance. These dual benefits highlight the computational advantages of inverse effectiveness principles in multimodal integration systems.

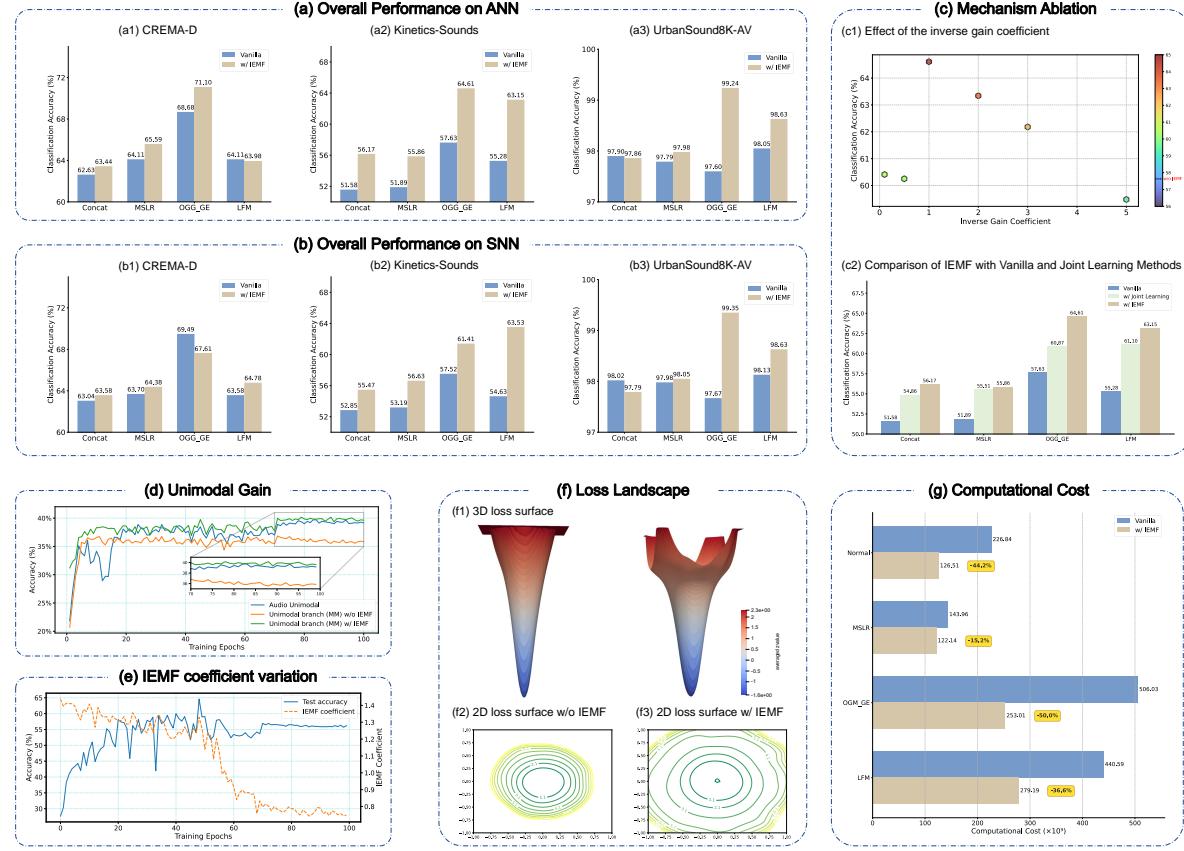
Overall, our contributions can be categorized into the following three points:

- We introduce the inverse effectiveness mechanism into multimodal fusion in deep neural networks for the first time, proposing an inverse effectiveness driven multimodal fusion strategy. This strategy adjusts the parameter update intensity of fusion modules in real-time, enabling the model to enhance its ability to extract information from other modalities when a single modality signal is weak, thereby improving information compensation effects and fusion efficiency.
- We evaluate our proposed method on two different architectures: Artificial Neural Networks (ANN) and Spiking Neural Networks (SNN). Experimental results demonstrate that IEMF possesses good generality and can be effectively integrated with both types of networks, fully leveraging its mechanism advantages.
- We conduct systematic empirical studies on multiple standard datasets and tasks, including representative scenarios such as audiovisual speech recognition, audiovisual continual learning, and audiovisual question answering. Experimental results show that our proposed method exhibits stronger perceptual capabilities under various complex conditions. Particularly worth emphasizing is that, as a mechanism, IEMF can seamlessly integrate with various existing state-of-the-art methods and further enhance their performance.

## 2. Results

### IEMF's generalizability across network architectures

A key advantage of IEMF lies in its strong universality across different neural network architectures. To evaluate this characteristic, we integrated IEMF into both Artificial Neural Networks (ANN) and Spiking Neural Networks



**Figure 2. Comprehensive evaluation of the proposed inverse effectiveness driven multimodal fusion (IEMF).** (a) Overall performance on ANNs. Bar charts compare the vanilla method (blue) with the method augmented by IEMF (khaki) on three audiovisual classification benchmarks—CREMA-D (a1), Kinetics-Sounds (a2) and UrbanSound8K-AV (a3)—under four representative fusion schemes (Concat, MSLR, OGM\_GE and LFM). (b) Overall performance on SNNs. Same layout as (a) but using spiking neural networks, demonstrating that IEMF consistently boosts accuracy across network paradigms and datasets (b1–b3). (c) Mechanism ablation on the Kinetics-Sounds dataset. (c1) Effect of the inverse gain coefficient  $\gamma$ : the baseline without IEMF (grey dashed line on the colour bar) scores below every IEMF setting; model accuracy peaks at  $\gamma=1$ . (c2) Removing the IEMF term (“Joint Learning only”) leads to a clear performance drop, highlighting the essential role of the inverse effectiveness multimodal fusion component. (d) Unimodal gain analysis. Test accuracy for (i) a unimodal audio model trained alone (“Unimodal”), (ii) the audio branch extracted from a multimodal model without IEMF (“Unimodal branch (MM) w/o IEMF”), and (iii) the audio branch with IEMF (“Unimodal branch (MM) w/ IEMF”). IEMF yields a persistent relative gain for the unimodal branch. (e) Dynamics of the IEMF coefficient. Evolution of the learnt IEMF coefficient  $\xi$  (dashed orange, right y-axis) alongside the test accuracy (solid blue, left y-axis) during training. At the early stage, the value of  $\xi$  is large to accelerate the multimodal integration, while as the network converges,  $\xi$  falls back and remains stable to maintain the fusion stability. (f) Loss landscape visualization. (f1) 3D loss surface comparison: vanilla method (left) versus IEMF-enhanced method (right). (f2-f3) 2D contour plots: without IEMF (f2) versus with IEMF (f3). The IEMF method leads to broader and flatter minima. (g) Computational cost analysis. Comparison of computational cost between standard models (blue) and IEMF-enhanced models (khaki). IEMF significantly reduces computational costs across all fusion methods, with reductions ranging from 15.2% to 50.0% (highlighted in yellow percentages).

(SNN), which represent distinctly different information processing paradigms. As shown in Fig. 2(a) and 2(b), IEMF consistently improves performance across multiple audio-visual classification benchmark tasks, regardless of the underlying network type. For example, with ANNs, IEMF increased classification accuracy from 51.58% to 56.17% under the Concat fusion on the Kinetics-Sounds dataset; similarly with SNNs, IEMF brought stable performance gains,

improving model accuracy from 52.85% to 55.47%, verifying the robustness and flexibility of our proposed method.

This cross-architecture robustness is particularly important for practical applications, as real-world systems often employ heterogeneous network models due to hardware resource constraints, power limitations, or real-time processing requirements. Notably, IEMF achieved significant benefits even in scenarios where traditional multimodal fusion

methods are limited by the sparsity and event-driven characteristics of spiking neural networks. For instance, when using the LFM fusion method on the Kinetics-Sounds dataset, the original SNN accuracy was 54.63%, slightly lower than the ANN's 55.28%; however, after introducing IEMF, the SNN classification accuracy surpassed the ANN, reaching 63.53% compared to the ANN's 63.15% with IEMF, as seen in Fig. 2(b2) and Fig. 2(a2). These results indicate that IEMF is not limited to traditional architectures but provides a universally applicable mechanism for improving multimodal fusion across various neural network architectures.

### IEMF improved the model performance on audio visual classification task

We systematically validated the effectiveness of the inverse effectiveness driven multimodal fusion (IEMF) mechanism in audio-visual classification tasks. We evaluated performance differences between baseline models and IEMF-enhanced models across three representative datasets: CREMA-D [9], Kinetics-Sounds [2], and UrbanSound8K-AV [16], using four mainstream fusion strategies: Concatenation Fusion (Concat), Modality-Specific Learning Rates (MSLR) [51], On-the-fly Gradient Modulation with Generalization Enhancement (OGM\\_GE) [37], and Learning Facilitator for Modality gap (LFM) [50]. As shown in Fig. 2(a1–a3), IEMF demonstrates consistent performance improvements across all fusion schemes and datasets.

Specifically, taking IEMF's enhancement to the MSLR method across datasets as an example, on the CREMA-D dataset, the baseline model achieved 64.11% accuracy using MSLR, which improved to 65.59% after introducing IEMF, yielding a 1.48% gain. On the more challenging Kinetics-Sounds dataset, baseline accuracy was 51.89%, while the IEMF-enhanced model reached 55.86%, a 3.97% improvement. Even on the UrbanSound8K-AV dataset where the baseline model already achieved high accuracy of 97.79%, IEMF further improved it to 97.98%. Though limited in magnitude, this improvement remains practically significant given the already high performance level.

It should be noted that in some cases, performance gains after introducing IEMF were relatively small, and in isolated instances even showed slight decreases (e.g., UrbanSound8K-AV dataset with Concat fusion strategy, Fig. 2(a3)). This phenomenon can primarily be attributed to: when baseline models already optimally leverage complementary audio-visual information in clear, low-noise environments, the existing modal contribution ratios are already near-optimal, naturally diminishing the benefits of dynamic adjustment and occasionally introducing slight perturbations due to additional modeling freedom. Therefore, IEMF's performance improvement potential is relatively limited in high-baseline, low-interference environments; whereas in environments with fluctuating modal sig-

nal quality or noise interference, IEMF's adaptive regulation mechanism demonstrates more significant advantages.

Looking at overall trends, IEMF consistently improves performance across datasets and fusion strategies, confirming its effectiveness in enhancing multimodal fusion efficiency. Unlike traditional fusion methods, IEMF dynamically adjusts modal fusion module weights based on each modality's relative strength. When information in one modality (e.g., audio) decreases due to noise or distortion, IEMF promotes greater information compensation from the fusion module, improving overall perceptual accuracy and model robustness. This dynamic adaptive mechanism significantly enhances model robustness when facing input quality fluctuations and environmental uncertainties.

To further validate IEMF's effectiveness, we conducted mechanism ablation experiments (Fig. 2(c1)-(c2)). In (c1), we analyzed classification accuracy changes under different inverse gain coefficient  $\gamma$  settings. Results show appropriate inverse gain effectively improves model performance, with optimal accuracy at  $\gamma = 1$ , indicating IEMF effectively balances unimodal and multimodal fusion signal contributions to maximize dynamic compensation. With larger coefficients (e.g.,  $\gamma = 5$ ), accuracy decreases, likely due to training instability from excessive modulation intensity. Conversely, without IEMF (w/o IEMF baseline in Fig. 2(c1)), classification accuracy is notably lower than all inverse gain coefficient settings, further validating the crucial role of IEMF in enhancing multimodal fusion. In (c2), we compared baseline models (Vanilla), models with joint learning strategy, and IEMF-enhanced models. We specifically included joint learning comparison to systematically evaluate IEMF's effectiveness. Joint learning adds independent classification heads for each modality without introducing new modalities, enhancing unimodal feature discriminability. IEMF dynamically modulates fusion module updates based on unimodal-fusion signal strength relationships for more refined compensation—mechanistically different approaches. Results show that while joint learning provides performance improvements, IEMF further enhances model performance, validating IEMF's superior generalizability through dynamic fusion module adjustment in existing multimodal learning frameworks.

We further evaluated how multimodal learning affects performance of weaker modality branches (audio) on Kinetics-Sounds using OGM\\_GE fusion (Fig. 2(d)). After multimodal training, we fine-tuned the audio branch to analyze fusion effects on unimodal perception. Results show traditional fusion methods (orange curve) lead to overfitting, limiting performance gains and even underperforming independently trained unimodal branches (blue curve). This suggests modal interference in conventional fusion degrades unimodal feature quality and perception. In contrast, with IEMF (green curve), the unimodal branch maintains

higher, more stable accuracy throughout training with significant early performance advantages. This confirms IEMF not only optimizes multimodal fusion but effectively mitigates modal interference, promoting better unimodal feature learning and generalization.

Figure 2(e) shows the test-set evolution of IEMF dynamic coefficients  $\xi$  and classification accuracy across training epochs. IEMF adapts fusion module behavior based on fusion effectiveness: during early training, fusion benefits are greater, keeping dynamic coefficient  $\xi$  high to maximize multimodal advantages; as unimodal features mature and fusion advantages diminish,  $\xi$  naturally decreases, reflecting reduced fusion dependency and helping maintain module stability for consistent test performance.

To validate the generalization properties of our proposed method, we visualized the loss landscapes of models with and without IEMF as shown in Fig. 2(f). The 3D loss surface visualization illustrated in Fig. 2(f1) reveals significant topological differences: the baseline method exhibits a sharper, cone-like minimum, while the IEMF-enhanced model displays a broader, more gradual basin structure. This distinction is further emphasized in the 2D contour plots depicted in Fig. 2(f2-f3): without IEMF, the contours form elongated elliptical patterns, indicating inconsistent curvature across different parameter directions; with IEMF, contours appear more circular and uniformly spaced, confirming a significantly flatter minimum region. These observations closely align with our subsequent theoretical analysis presented later in this paper, which demonstrates that IEMF directs the optimization process toward flatter regions of the loss landscape, a characteristic directly associated with the improved generalization performance observed in our experimental results.

Beyond performance improvements, we analyzed IEMF’s impact on computational cost as shown in Fig. 2(g). Our evaluation employed a comprehensive computational cost metric that balances both convergence speed and per-epoch complexity, providing a more holistic assessment of algorithmic efficiency. Across all fusion methods, IEMF consistently reduces computational costs by significant margins. The computational savings range from 15.2% for MSLR to 50.0% for OGM.GE, with Normal and LFM configurations showing reductions of 44.2% and 36.6%, respectively. These substantial improvements stem from IEMF’s ability to achieve faster convergence while maintaining reasonable per-epoch complexity. By dynamically modulating fusion behavior based on modality contributions, IEMF effectively reduces the total computational budget required to reach optimal performance. Importantly, these efficiency gains occur concurrently with the performance enhancements reported earlier, demonstrating that IEMF not only improves model accuracy but also significantly optimizes computational resource utilization—a crit-

ical advantage for resource-constrained multimodal applications in real-world environments.

In summary, IEMF demonstrates consistent performance improvements across datasets and fusion strategies. Systematic experiments validate its effectiveness in dynamically regulating fusion, mitigating modal interference, enhancing unimodal learning, and improving robustness, providing an efficient and well-generalizing fusion strategy for multimodal perception tasks.

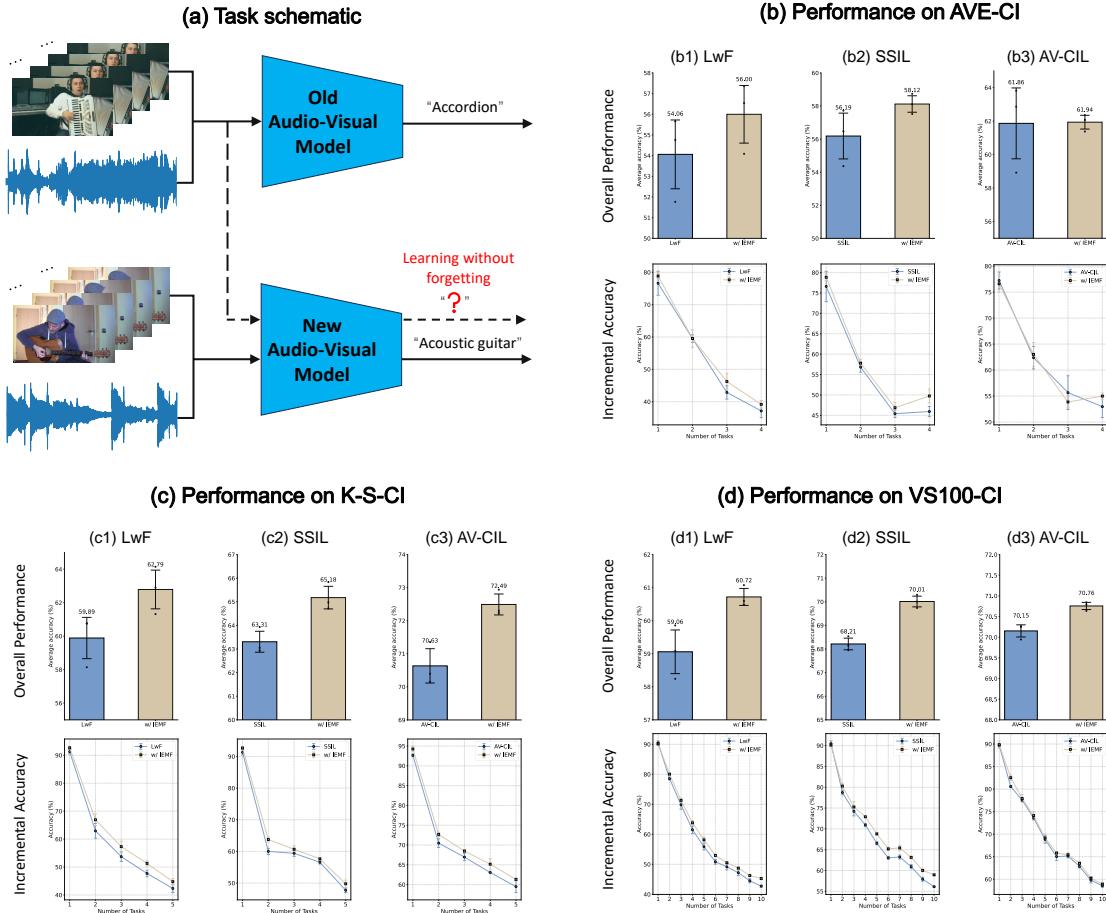
### IEMF improved the model performance on audio visual continual learning

To evaluate the effectiveness of IEMF in more challenging scenarios, we further examined its performance in audio-visual continual learning tasks. In such tasks, models need to learn new categories continuously while preserving recognition capabilities for previously learned categories as much as possible, thereby avoiding catastrophic forgetting, as shown in Fig. 3(a). We selected three representative class-incremental learning baseline methods for comparison: LwF [27], SSIL [1], and AV-CIL [38], and evaluated them on three audio-visual continual learning datasets: AVE-CI, K-S-CI, and VS100-CI [38]. The experimental results are shown in Fig. 3(b-d).

After introducing the IEMF method, the models achieved stable accuracy improvements across all datasets. On AVE-CI, LwF increased from 54.06 % to 56.00 % (+1.94 %), SSIL improved from 56.19 % to 58.12 % (+1.93 %), and AV-CIL slightly increased from 61.86 % to 61.94 % (+0.08 %). In K-S-CI, which features more cross-modal noise, LwF rose from 59.89 % to 62.79 % (+2.90 %), SSIL improved from 63.31 % to 65.18 % (+1.87 %), and AV-CIL increased from 70.63 % to 72.49 % (+1.86 %). For the largest scale dataset VS100-CI, LwF improved from 59.06 % to 60.72 % (+1.66 %), SSIL from 68.21 % to 70.01 % (+1.80 %), and AV-CIL from 70.15 % to 70.76 % (+0.61 %). All nine comparisons showed positive gains, with an average improvement of approximately 1.63 percentage points, highlighting the consistent effectiveness of IEMF.

In Fig. 3(b-d), the line graphs in the bottom row of each subfigure show the average accuracy changes after each continuous task. Notably, compared to baseline models, the accuracy decline curves of IEMF models are significantly more gradual. This indicates that IEMF enhances the model’s ability to retain existing knowledge during cross-task knowledge transfer while effectively integrating information about new categories, thereby significantly mitigating catastrophic forgetting.

To further understand the internal mechanisms behind IEMF’s performance improvements, we analyzed its fusion dynamic behavior during training. IEMF does not explicitly introduce learnable parameters bound to specific tasks, but rather adaptively regulates the update dynamics of the



**Figure 3. Inverse effectiveness driven multimodal fusion boosts audio visual continual learning.** (a) **Task schematic.** A single audiovisual model is incrementally updated as new classes arrive; the goal is to absorb the new knowledge while preserving performance on previously learned classes—achieving “learning without forgetting”. (b) **Results on AVE-CI, (c) K-S-CI and (d) VS100-CI.** For three representative class incremental learning baselines—LwF, SSIL and AV-CIL—we compare the vanilla method (blue) with the method augmented by IEMF (khaki). Each sub-panel is split into top and bottom: the top bar chart reports the *overall performance* (mean accuracy across all tasks, error bars denote one standard deviation), while the bottom line plot traces the *incremental accuracy* after each successive task. Across all datasets and baselines, IEMF consistently increases mean accuracy and yields a flatter accuracy-decay curve, indicating the better knowledge transfer.

fusion module based on changes in the effectiveness of unimodal and multimodal signals, thereby implicitly adapting to modal variations across different task stages during continuous learning. Through the training process guided by the inverse effectiveness principle, the model can naturally adapt to changes in modal reliability during weight updates, thus reducing over-reliance on a single modality when perceptual conditions fluctuate. Benefiting from this adaptive optimization strategy, IEMF not only improves the average accuracy across all tasks but also maintains a smoother performance degradation curve, preserving high overall performance and cross-task knowledge coherence even as new tasks are continuously introduced.

### IEMF improved the model performance on audio visual question answering

We further evaluated the effectiveness of IEMF in audio visual question answering (AVQA) tasks. In this task, models must answer text questions based on synchronized audio and video inputs, demanding higher capabilities for deep integration of multimodal information. As shown in Fig. 4, radar charts 4(a1) and 4(b1) compare the classification accuracy of baseline models versus models with IEMF, and ST-AVQA [26] models versus models with IEMF, across different question types (audio-only questions, visual-only questions, and audio-visual combined questions).

Comparing the radar charts of original models and models with IEMF, we can observe that IEMF improved an-

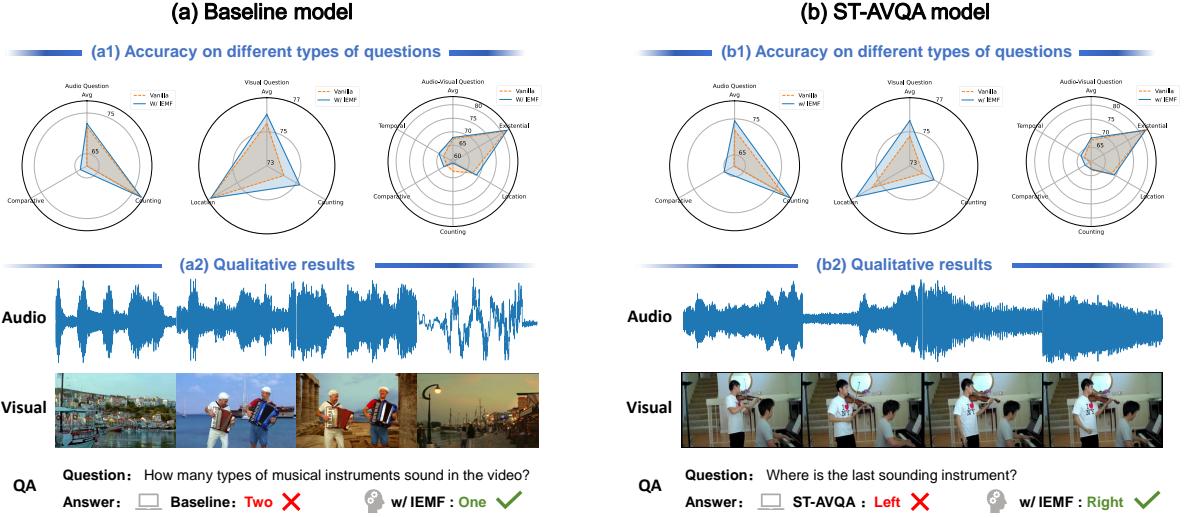


Figure 4. **Quantitative and qualitative impact of IEMF on audio visual question answering task.** **(a)** Baseline model. The three radar charts (top) report accuracy on audio-only, visual-only, and audio-visual questions, respectively. Orange = vanilla, Blue = w/ IEMF. The bottom row shows a representative sample—waveform, video frames, and question/answer—where the vanilla fusion miscounts the instruments (“Two”), whereas IEMF answers correctly (“One”). **(b)** ST-AVQA model. IEMF again enlarges the radar area for every question type and corrects the localization query in the illustrated example (vanilla: “Left”; w/ IEMF: “Right”). Across both models, the blue polygons consistently enclose the orange ones, confirming that the inverse effectiveness driven multimodal fusion mechanism improves all question categories while providing intuitive per-sample gains.

swer accuracy across all question types. Taking the ST-AVQA model and its IEMF-enhanced version as an example (Fig. 4(b)), for audio-only questions, the original ST-AVQA model achieved an average accuracy of 71.90%, while the IEMF model improved to 74.49%, an increase of 2.59%. Similarly, for visual-only questions, the baseline accuracy was 74.74%, while the IEMF-enhanced model reached 75.65%, an improvement of 0.91%. For audio-visual questions, the vanilla model’s average accuracy was 67.61%, while the IEMF model achieved 68.33%, an improvement of 0.72%. To verify IEMF’s performance on fine-grained questions, we specifically analyzed its effectiveness in tasks requiring precise localization classification. As shown in Fig. 4(b2), the original ST-AVQA model incorrectly predicted “left” side when answering “the position of the last sounding instrument”, while the model with IEMF correctly located it as “right” side. This demonstrates that IEMF-enhanced models possess stronger fine-grained discrimination capabilities in complex cross-modal reasoning tasks, improving the integration efficiency of multimodal cues.

This improvement further validates IEMF’s crucial advantages in handling noisy interference or incomplete input information scenarios. By dynamically adjusting the update rate of the fusion module during training based on the strength of unimodal and multimodal fusion signals, IEMF guides the model to learn strategies that can more robustly integrate different modal information when modal signal

strengths are uneven or information is contradictory. In contrast, models without IEMF are more prone to judgment biases when facing modal conflicts or input uncertainties, leading to incorrect answers or overall performance degradation. Overall, these results highlight IEMF’s important role in enhancing multimodal understanding and reasoning.

### 3. Discussion

#### Biological insights into multimodal integration

Despite significant advances in multimodal fusion, many key biological principles have not been fully explored and applied in artificial intelligence systems, which could further enhance the robustness and adaptability of multimodal systems. In this study, we propose a brain-inspired deep neural network multimodal integration method based on the inverse effectiveness mechanism observed in biological systems, providing verifiable theoretical foundations and methodological support for multimodal information integration. Specifically, we propose an inverse effectiveness driven multimodal fusion (IEMF) method that dynamically adjusts the weights of modal fusion modules based on the relationship between the strength of single-modal cues and the signal strength after modal fusion. Our approach systematically considers the complementary interactions between modalities, significantly improving not only the performance and generalization capabilities of multimodal systems but also their computational efficiency. This dual ad-

vantage—enhanced robustness coupled with reduced computational costs—offers insights into why inverse effectiveness might have evolved as a critical mechanism in biological systems, where both perceptual reliability and metabolic efficiency are under evolutionary pressure.

In the IEMF framework, we prioritize inverse effectiveness in the network training process by introducing an inverse effectiveness gain coefficient that applies gradient regulation to fusion weights, naturally forming an internal bias of “weak modality-high gain, strong modality-low gain” during the learning phase. This strategy aligns with the physiological mechanism of cross-modal experience shaping plasticity in early neural development in biological organisms, where newborn individuals initially lack multisensory integration abilities. These abilities are not innate but develop through continuous shaping of neural circuits via early cross-modal experiences, adapting to the environment and optimizing multimodal integration performance [46]. Furthermore, we apply inverse effectiveness driven fusion strategies uniformly across all input channels of the model, consistent with the view in [39] that inverse effectiveness manifests not only under degraded stimulus conditions but also with clear stimuli.

This research confirms the long-standing intuition that how modalities are combined is as important as how many are combined. By introducing inverse effectiveness rules from cortical circuits into gradient-based optimization learning systems, we achieved: (i) effective generalization in both Artificial Neural Networks (ANN) and Spiking Neural Networks (SNN) models; (ii) significant performance improvements in audio-visual classification, audio-visual continual learning and audio-visual question answering tasks; and (iii) substantial computational efficiency gains with up to 50% reduction in computational costs across diverse fusion methods. Systematic experiments demonstrate that after integrating the IEMF mechanism into existing multimodal methods, models achieved performance superior to original state-of-the-art techniques across various multimodal tasks, further indicating that introducing bio-inspired mechanisms can effectively improve the efficiency of multimodal integration, expanding its potential in artificial intelligence applications. These findings not only highlight the advantages of incorporating biological principles into machine learning models but also provide new directions for future research in neuromorphic computing and multisensory integration.

### Other biological mechanisms of multimodal integration

It is worth discussing that while this research emphasizes the importance of inverse effectiveness in multimodal integration, it is worth noting that there are two other equally important principles in multimodal biological perception processes: temporal congruence and spatial congruence.

These mechanisms are particularly important in dynamic multimodal integration. Temporal congruence refers to visual and auditory inputs maintaining coordination in time, thereby optimizing perceptual and decision-making performance. Experimental studies [17, 43] demonstrate that when visual and auditory stimuli are presented in close synchronization within a 0-200 millisecond time window, they significantly enhance the accuracy and reaction speed of perceptual judgments. In contrast, temporal asynchrony leads to decreased activation intensity in relevant brain regions, weakening the integration effect. Spatial congruence refers to different sensory modalities maintaining consistency or proximity in spatial location, thereby enhancing the joint representation of cross-modal signals. Research has found that multisensory neurons (such as neurons in the superior colliculus) exhibit integration enhancement effects only when audiovisual stimuli originate from the same or adjacent spatial locations; otherwise, integration may be inhibited or show no integration response [44]. For example, in real-world scene understanding, object recognition, and tracking tasks, accurately matching sound sources with corresponding visual objects is key to successful multimodal perception. Temporal and spatial congruence are crucial for accurate multimodal integration.

Although temporal and spatial congruence are indispensable in biological perception, given that the tasks selected in this study inherently possess strong input synchronization characteristics (i.e., dual-modal inputs from the same source at the same moment), we did not explicitly model these mechanisms in the current work. Specifically, the sensory inputs in this study’s tasks naturally possess synchronization and correspondence relationships; therefore, these congruence factors have already been implicitly considered in the multimodal fusion process. Looking forward, further research could explore how to explicitly incorporate temporal and spatial congruence into the IEMF framework by introducing asynchronous, spatially disparate multimodal input samples, thereby training models to effectively integrate under more complex temporal and spatial variation conditions, further advancing biologically inspired multimodal learning systems toward broader application domains.

## 4. Materials and methods

### Neuron models in ANNs and SNNs

In neural networks, the information flow is governed by the dynamics of neuronal activation. In this work, we adopt two distinct neuron models: the continuous artificial neurons used in artificial neural networks (ANNs) and the spike-based neurons in spiking neural networks (SNNs).

In ANNs, information is processed continuously. Each neuron computes a weighted sum of its inputs and applies a nonlinear activation function to produce its output:  $y = f(Wx+b)$ , where  $x$  is the input vector,  $W$  is the weight

matrix,  $b$  is the bias term, and  $f(\cdot)$  denotes a nonlinear function such as ReLU [33] or sigmoid.

In contrast, SNNs more closely mimic biological neurons by communicating via discrete spike events. We employ the widely used Leaky Integrate-and-Fire (LIF) model [11] to capture the membrane potential dynamics. Upon receiving a synaptic input current  $I(t)$ , the membrane potential  $U(t)$  accumulates over time. When  $U(t)$  crosses a threshold  $U_{\text{th}}$ , a spike is emitted and the potential is reset to the resting value  $U_{\text{rest}}$ . The continuous-time dynamics of the LIF neuron are given by:

$$\tau_m \frac{dU(t)}{dt} = -(U(t) - U_L) - \frac{g_{E|I}}{g_L}(U(t) - U_{E|I}) + \frac{I_s}{g_L}, \quad (1)$$

where  $\tau_m = C_m/g_L$  is the membrane time constant,  $C_m$  is membrane capacitance,  $g_L$  is the leak conductance, and  $g_{E|I}$  and  $U_{E|I}$  denote the conductance and reversal potentials for excitatory or inhibitory synapses.  $I_s$  is the synaptic input current. To simplify the formulation, we aggregate the synaptic terms into an effective input current:  $RI(t) \triangleq -\frac{g_{E|I}}{g_L}(U(t) - U_{E|I}) + \frac{I_s}{g_L}$ , reducing the membrane potential equation to:

$$\tau_m \frac{dU(t)}{dt} = -(U(t) - U_{\text{rest}}) + RI(t). \quad (2)$$

For numerical simulation, we set  $U_{\text{rest}} = 0$  and discretize the above dynamics. To clearly distinguish continuous and discrete states, we denote membrane potential as  $\mathbf{u}^t$  at discrete time step  $t$ . The complete discrete-time update of the membrane potential and spike generation at layer  $l$  is:

$$\begin{cases} \mathbf{u}_{\text{pre}}^{t,l} = \tau \mathbf{u}^{t-1,l} + \mathbf{W}^l \mathbf{s}^{t,l-1}, & (\text{accumulation}) \\ \mathbf{s}^{t,l} = H(\mathbf{u}_{\text{pre}}^{t,l} - \mathbf{u}_{\text{th}}), & (\text{spike firing}) \\ \mathbf{u}^{t,l} = \mathbf{u}_{\text{pre}}^{t,l} \cdot (1 - \mathbf{s}^{t,l}), & (\text{reset mechanism}) \end{cases} \quad (3)$$

where  $\mathbf{W}^l$  is the weight matrix from layer  $l-1$  to  $l$ ,  $\mathbf{s}^{t,l-1}$  denotes the spike train from the previous layer at time  $t$ ,  $\tau = 1 - \frac{1}{\tau_m}$  is the leak factor controlling the temporal decay of the membrane potential, and  $H(\cdot)$  is the Heaviside step function used to generate binary spike outputs.

### Multimodal integration formulation

We denote a multimodal input as  $\mathbf{x} = (\mathbf{x}^a, \mathbf{x}^v)$ , where  $\mathbf{x}^a \in \mathcal{X}^a$  and  $\mathbf{x}^v \in \mathcal{X}^v$  represent inputs from two modalities (e.g., audio and visual). These are independently processed by two encoders  $\varphi^a(\cdot; \theta^a)$  and  $\varphi^v(\cdot; \theta^v)$  to obtain modality-specific latent representations:

$$\mathbf{z}^a = \varphi^a(\mathbf{x}^a; \theta^a), \quad \mathbf{z}^v = \varphi^v(\mathbf{x}^v; \theta^v), \quad (4)$$

where  $\theta^a$  and  $\theta^v$  represent the trainable parameters of the audio and visual encoders, respectively. The extracted features are fused using a general fusion operator  $\mathcal{F}(\cdot, \cdot)$ , i.e.,

$\mathbf{z}^{\text{av}} = \mathcal{F}(\mathbf{z}^a, \mathbf{z}^v)$ , followed by a classifier  $h(\cdot; \theta^h)$  that maps the fused audio visual features to a prediction  $\hat{\mathbf{y}} = h(\mathbf{z}^{\text{av}}; \theta^h)$ , where  $\theta^h$  denotes the classifier parameters.

In the widely adopted vanilla fusion strategy, such as feature concatenation, the fusion operator takes the form:

$$\mathcal{F}(\mathbf{z}^a, \mathbf{z}^v) = \mathbf{W}^f [\mathbf{z}^a; \mathbf{z}^v] + \mathbf{b}^f, \quad (5)$$

where  $[\cdot; \cdot]$  denotes the concatenation operation,  $\mathbf{W}^f \in \mathbb{R}^{M \times (d_a+d_v)}$  and  $\mathbf{b}^f \in \mathbb{R}^M$  are the parameters of the fusion layer, and  $M$  is the number of output classes.

The multimodal learning goal is to train a multimodal model  $f_{\theta} : \mathcal{X}^a \times \mathcal{X}^v \rightarrow \mathcal{Y}$ , where the learnable parameters  $\theta = \{\theta^a, \theta^v, \theta^h\}$ , that minimizes the empirical risk over a dataset  $\mathcal{D} = (\mathbf{x}_i^a, \mathbf{x}_i^v, y_i)_{i=1}^N$ . The training objective is:

$$\arg \min_{\theta} \mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\text{ce}}(f_{\theta}(\mathbf{x}_i^a, \mathbf{x}_i^v), y_i), \quad (6)$$

where  $\mathcal{L}_{\text{ce}}(\cdot)$  denotes the cross-entropy loss function.

### Inverse effectiveness driven multimodal fusion

Most previous studies have focused on designing sophisticated fusion blocks, refining the joint representation  $\mathbf{z}^{\text{av}}$  to enhance cross-modal interactions, yet they rarely consider how the relative informativeness of each unimodal stream relative to the fused output should guide the fusion module. In this work, we draw inspiration from the principle of inverse effectiveness. Instead of employing static cross-modal integration, we adaptively adjust the update rate of the fusion module's weights by contrasting the estimated information content of each unimodal branch with that of the integrated multimodal representation.

First, for each sample  $i$  in a mini-batch  $\mathcal{B}_t$ , we evaluate the per-sample modal information content  $c_i$  according to the following equation:

$$\begin{cases} c_i^a = [\mathbf{p}_i^a]_{y_i}, & \mathbf{p}_i^a = \pi(\mathbf{W}_t^a \cdot \mathbf{z}^a + \mathbf{b}_t^a), \\ c_i^v = [\mathbf{p}_i^v]_{y_i}, & \mathbf{p}_i^v = \pi(\mathbf{W}_t^v \cdot \mathbf{z}^v + \mathbf{b}_t^v), \end{cases} \quad (7)$$

where  $\mathbf{W}_t^{a/v}$ ,  $\mathbf{b}_t^{a/v}$  are the parameters of the audio/visual modal classification heads, respectively, and  $[\mathbf{p}]_{y_i}$  picks the probability assigned to the ground-truth label  $y_i$ .  $\pi$  is a normalization function, and here we choose the softmax function. The informativeness of the multimodal output is estimated in the same way:

$$c_i^{\text{av}} = [\mathbf{p}_i^{\text{av}}]_{y_i}, \quad \mathbf{p}_i^{\text{av}} = \pi(\mathbf{W}_t^{\text{av}} \cdot \mathbf{z}^{\text{av}} + \mathbf{b}_t^{\text{av}}), \quad (8)$$

Next, we average the evaluated values in (7) and (8) to obtain the batch-level modality-strength scores.

$$S_t^{a-v} = \frac{1}{2|\mathcal{B}_t|} \sum_{i \in \mathcal{B}_t} (c_i^a + c_i^v), \quad S_t^{\text{av}} = \frac{1}{|\mathcal{B}_t|} \sum_{i \in \mathcal{B}_t} c_i^{\text{av}}. \quad (9)$$

Following the biological observation that fusion should dominate, and provides the greatest benefit when unimodal evidence is weak, we define a inverse effectiveness driven multimodal fusion coefficient  $\xi_t$ . The IEMF coefficient quantifies how effective the unimodal branches are relative to the fused output. We map  $\xi_t$  to a bounded value:

$$\xi_t = \gamma \cdot \left( 1 + \kappa \left( 1 - \frac{S_t^{a-v}}{S_t^{av}} \right) \right), \quad \gamma > 0, \quad (10)$$

where  $\gamma$  is the inverse gain coefficient controlling the overall magnitude of fusion modulation, and  $\kappa(\cdot)$  denotes a generic bounded gating function, in this work we instantiate  $\kappa$  with the hyperbolic tangent, i.e.,  $\kappa(\cdot) = \tanh(\cdot)$ , owing to its smooth and symmetric saturation properties. Because  $\kappa(\cdot) \in (-1, 1)$ , Eq. (10) confines the fusion coefficient to the interval  $\xi_t \in (0, 2\gamma)$ .

IEMF coefficient  $\xi_t$  magnitude varies intuitively with the strength ratio between unimodal and multimodal evidence:

- *Weak unimodal evidence* ( $S_t^{a-v} \ll S_t^{av}$ ). The term  $1 - S_t^{a-v}/S_t^{av}$  is positive and thus  $\xi_t$  approaches its upper bound. A larger  $\xi_t$  amplifies the fusion gradient, encouraging the model to rely more heavily on cross-modal integration.
- *Strong unimodal evidence* ( $S_t^{a-v} \gtrsim S_t^{av}$ ). As the ratio nears or exceeds 1, the inner term becomes non-positive;  $1 + \kappa(\cdot)$  decreases, pulling  $\xi_t$  toward its lower limit 0. When the unimodal score predominates, a smaller  $\xi_t$  attenuates the fusion update, preserving the integrity of an already robust unimodal pathway.

In summary, a large  $\xi_t$  corresponds to weak unimodal cues and triggers a stronger adjustment of the fusion weights, whereas a small  $\xi_t$  indicates confident unimodal predictions and results in a milder fusion update.

The inverse effectiveness driven multimodal fusion coefficient  $\xi_t$  is applied only to the fusion parameters and the unimodal branches are not affected. Concretely,

$$\mathbf{W}_{t+1}^f = \mathbf{W}_t^f - \eta \xi_t \nabla_{\mathbf{W}^f} \mathcal{L}(\mathbf{W}_t^f), \quad (11)$$

where  $\eta$  is the learning rate and  $\nabla_{\mathbf{W}^f} \mathcal{L}(\mathbf{W}_t^f)$  is the raw fusion gradient. Because  $0 < \xi_t < 2\gamma$ , this scaling never reverses the descent direction, thereby maintaining optimisation stability. Together, Eqs. (7)–(11), inspired by the inverse effectiveness principle, are characterized as: the fusion pathway receives a larger update when unimodal evidence is weak and a smaller one when unimodal confidence is high. This self-balancing rule enhances robustness under noise and improves the model’s generalization across diverse input conditions.

### Theoretical analysis of inverse effectiveness driven multimodal fusion strategy

We prove that the inverse-effectiveness coefficient  $\xi_t$  used by IEMF reduces the expected step size more in high-curvature directions, ensuring reliable convergence to local minima while maintaining optimization stability throughout the training process.

**Assumption 1.** *The loss  $\mathcal{L}(\mathbf{W}^f)$  is twice continuously differentiable and there exist constants  $\beta, \rho > 0$  such that for all  $\mathbf{u}, \mathbf{v}$ ,  $\|\nabla \mathcal{L}(\mathbf{u}) - \nabla \mathcal{L}(\mathbf{v})\| \leq \beta \|\mathbf{u} - \mathbf{v}\|$ ,  $\|\mathbf{H}(\mathbf{u}) - \mathbf{H}(\mathbf{v})\| \leq \rho \|\mathbf{u} - \mathbf{v}\|$ . This means that the gradient and the Hessian are  $\beta$ -Smoothness and  $\rho$ -Lipschitz respectively.*

**Theorem 1** (Convergence properties of IEMF). *Assume 1. Let  $\mathbf{W}^{f*}$  be a local minimizer and write  $\mathbf{H}^* = \mathbf{H}(\mathbf{W}^{f*})$  with eigen-pairs  $\{(\lambda_i, \mathbf{e}_i)\}_{i=1}^d$ ,  $0 < \lambda_1 \leq \dots \leq \lambda_d$ . The IEMF updates the fusion module parameters according to the equation (11). Define the deviation  $\Delta_t := \mathbf{W}_t^f - \mathbf{W}^{f*} = \sum_{i=1}^d \alpha_i^t \mathbf{e}_i$ . Choose a radius  $r > 0$  such that*

$$\frac{\rho}{2} r < \lambda_1 \quad \text{and} \quad \|\Delta_t\| \leq r \quad \forall t.$$

*Then,*

$$\mathbb{E}[\eta \xi_t \lambda_i] \begin{cases} < \eta \lambda_i, & S_t^{a-v}/S_t^{av} > 1 \quad (\text{unimodal dominated}), \\ = \eta \lambda_i, & S_t^{a-v}/S_t^{av} = 1, \\ > \eta \lambda_i, & S_t^{a-v}/S_t^{av} < 1 \quad (\text{fusion dominated}). \end{cases}$$

*As a result, with IEMF, unimodal-dominated batches reduce sharp directions more than vanilla method, whereas fusion-dominated batches allow at most a two-fold increase in step size, thus preserving optimization convergence.*

*Proof.* For any  $\mathbf{W}^f$  satisfying  $\|\mathbf{W}^f - \mathbf{W}^{f*}\| \leq r$ , we use  $\nabla \mathcal{L}(\mathbf{W})$  denotes  $\nabla_{\mathbf{W}^f} \mathcal{L}(\mathbf{W}^f)$  and construct  $g(s) = \nabla_{\mathbf{W}^f} \mathcal{L}(\mathbf{W}^{f*} + s\Delta)$ ,  $s \in [0, 1]$ . The fundamental theorem of calculus gives  $\nabla \mathcal{L}(\mathbf{W}) = g(1) - g(0) = \int_0^1 \frac{d}{ds} g(s) ds$ . Then we have:

$$\nabla \mathcal{L}(\mathbf{W}) = \int_0^1 \mathbf{H}(\mathbf{W}^f + s\Delta) \Delta ds.$$

Adding and subtracting  $\mathbf{H}^* \Delta$  inside the integral yields:

$$\nabla \mathcal{L}(\mathbf{W}) = \mathbf{H}^* \Delta + \int_0^1 [\mathbf{H}(\mathbf{W}^{f*} + s\Delta) - \mathbf{H}^*] \Delta ds.$$

We define the remainder term as:

$$\mathbf{R}(\mathbf{W}) = \int_0^1 [\mathbf{H}(\mathbf{W}^{f*} + s\Delta) - \mathbf{H}^*] \Delta ds,$$

where the remainder  $\mathbf{R}(\mathbf{W})$  can be estimated for an upper bound with the help of the  $\rho$ -Lipschitz condition

$$\|\mathbf{R}(\mathbf{W})\| \leq \frac{\rho}{2} \|\mathbf{W} - \mathbf{W}^{f*}\|^2. \quad (12)$$

which allows us to express the gradient as:

$$\nabla \mathcal{L}(\mathbf{W}) = \mathbf{H}^*(\mathbf{W} - \mathbf{W}^{f*}) + \mathbf{R}(\mathbf{W}).$$

Insert the above equation into the update (11), and Take inner product with  $\mathbf{e}_i$  and use  $(\mathbf{H}^* \mathbf{e}_i)^\top = (\lambda_i \mathbf{e}_i)^\top$ , we have

$$\alpha_i^{t+1} = (1 - \eta \xi_t \lambda_i) \alpha_i^t - \eta \xi_t \mathbf{e}_i^\top \mathbf{R}(\mathbf{W}_t^f), \quad (13)$$

where  $\alpha_i^t = \mathbf{e}_i^\top \Delta_t$ . For the contraction argument we need the last term  $\eta \xi_t |\mathbf{e}_i^\top \mathbf{R}(\mathbf{W}_t^f)|$  to be strictly smaller than the linear part  $\lambda_i \|\Delta_t\|$  even in the flattest direction ( $\lambda_i = \lambda_1$ ). With Cauchy–Schwarz formula and (12), we have

$$\begin{aligned} |\mathbf{e}_i^\top \mathbf{R}(\mathbf{W}_t^f)| &\leq \|\mathbf{R}(\mathbf{W}_t^f)\|, \\ \|\mathbf{R}(\mathbf{W}_t^f)\| &\leq \frac{\rho}{2} \|\Delta_t\|^2. \end{aligned}$$

Applying the triangle inequality to Eq. (13) and substituting our derived bounds, we have

$$|\alpha_i^{t+1}| \leq |1 - \eta \xi_t \lambda_i| |\alpha_i^t| + \eta \xi_t \left( \frac{\rho}{2} \|\Delta_t\|^2 \right)$$

By imposing  $\frac{\rho}{2} r < \lambda_1$ , we establish a crucial inequality:

$$\eta \xi_t \left( \frac{\rho}{2} \|\Delta_t\|^2 \right) \leq \eta \xi_t \frac{\rho}{2} r \|\Delta_t\| < \eta \xi_t \lambda_1 \|\Delta_t\|.$$

This inequality demonstrates that the quadratic remainder term is always strictly dominated by the linear term for all eigendirections  $i$ , since  $\lambda_i \geq \lambda_1$  for all  $i$ . Consequently, the convergence behavior of each component  $\alpha_i^t$  is primarily determined by the multiplicative factor  $(1 - \eta \xi_t \lambda_i)$ , with the remainder term providing a bounded perturbation that does not disrupt the overall convergence pattern established by the linear term.  $\xi_t$  is computed by equation (10), with  $\gamma = 1, 0 < \xi_t < 2$ . Taking expectation of  $\eta \xi_t \lambda_i$  over the mini-batch three cases arise, yielding exactly the inequalities stated in the theorem.  $\square$

Our analysis reveals two key aspects of IEMF’s convergence properties, primarily manifested through the factor  $(1 - \eta \xi_t \lambda_i)$ , which precisely controls how quickly the error components  $\alpha_i^t$  (representing the projection of parameter error onto each eigenvector) contract toward zero. In unimodal-dominated batches ( $S_t^{a-v}/S_t^{av} > 1$  resulting in  $\xi_t < 1$ ), the contraction factor satisfies  $|1 - \eta \xi_t \lambda_i| < |1 - \eta \lambda_i|$ , meaning that high-curvature directions (larger  $\lambda_i$  values) contract faster than in vanilla method. Meanwhile, in fusion-dominated batches ( $S_t^{a-v}/S_t^{av} < 1$  resulting in  $\xi_t \in (1, 2)$ ), although step sizes may increase, they remain strictly bounded since  $\xi_t < 2$ , ensuring the algorithm’s global stability. This dual mechanism, which balances a preference for minima with strict step-size constraints, ensures that IEMF maintains reliable convergence properties while adaptively adjusting step sizes.

Empirical studies have demonstrated a strong correlation between flatter minima and improved generalization performance [15, 23]. While Theorem 1 establishes the convergence of IEMF under standard smoothness assumptions, a stronger theoretical link between the geometric properties of the solution and its generalization remains challenging to formally prove. Nevertheless, we provide the following heuristic justification based on landscape sharpness analysis, supported by experimental observations.

Consider the sharpness of the loss landscape at a parameter configuration  $\mathbf{W}^f$ , defined as

$$s(\mathbf{W}^f, \rho) := \max_{\|\epsilon\|_2 \leq \rho} \mathcal{L}(\mathbf{W}^f + \epsilon) - \mathcal{L}(\mathbf{W}^f), \quad (14)$$

which quantifies the sensitivity of the loss to local perturbations. While we do not provide a formal guarantee, our empirical evidence and directional analysis suggest the following heuristic conclusion:

**Heuristic 1** (IEMF reduces landscape sharpness). *Under dynamic training, IEMF adaptively reduces the step size in high-curvature directions. As a result, the sharpness*

$$\mathbb{E}[s(\mathbf{W}^f, \rho)] \lesssim \alpha \cdot s_{\text{vm}}(\mathbf{W}^f, \rho), \quad (15)$$

where  $s_{\text{vm}}$  is the sharpness observed under vanilla method and  $\alpha < 1$  is a factor that quantifies how much IEMF reduces the loss landscape’s sharpness through its adaptive modulation of optimization steps.

This suggests that IEMF biases the optimization trajectory toward flatter regions of the loss landscape, a property that empirically correlates with improved generalization.

## Experimental settings and training details

**Datasets.** We divided the dataset as specified in the original dataset. Audio-visual classification: CREMA-D [9], an audiovisual dataset containing six most common emotion categories for speech emotion recognition with total 7442 video clips. We randomly divided the dataset into a training and validation set, as well as a test set, with a ratio of 9:1; Kinetics-Sounds [2], contains 31 human action categories selected from the Kinetics dataset [22]. The dataset contains 17,366 10-second video clips, of which 1,472 are training and validation samples and 2,594 are test samples; UrbanSound8K-AV dataset [16], with 8732 audio-visual samples totaling 10 categories. Each sample consists of a color image and a 4-second audio signal. We randomly divided the dataset into training and test sets in the ratio of 7:3. Audio Visual Continual learning: We used the class-incremental audiovisual dataset introduced by [38], comprising three benchmark datasets: AVE-CI (4 tasks  $\times$  7 classes) consisting of 3,294 training samples, 391 validation samples, and 394 test samples, K-S-CI (5 tasks  $\times$  6 classes)

containing 19,220 training samples, 1,947 validation samples, and 1,958 test samples and VS100-CI (10 tasks  $\times$  10 classes) with 51,195 training samples, 5,000 validation samples, and 5,000 test samples. For audio-visual question answering, we used the official MUSIC-AVQA [26] split, which contains 32,087 training, 4,595 validation, and 9,185 test question-answer pairs.

**Data processing and network backbones.** All raw videos were first resampled to a uniform frame rate. According to different task settings, we randomly sampled 1, 3, or 16 frames from each video clip as the visual input. The corresponding audio input was transformed into log-Mel spectrograms, which were used as input to the audio branch. For audiovisual classification, we employed the ResNet-18 [18] architecture for both the visual and audio streams. To investigate architectural generality, we adapted this topology to a spiking neural network counterpart by replacing conventional activation functions with leaky integrate-and-fire (LIF) neurons; neuron-level hyper-parameters are detailed in Table S1. For audiovisual continual learning, we used VideoMAE [48] and AudioMAE [20] to extract video frames and audio features. For audio-visual question and answer, for vision we used pre-trained ResNet-18 model and for audio we used pre-trained VGGish [19] to extract visual and audio features respectively.

**Optimization details.** We trained models with stochastic gradient descent (SGD) [40] and a weight-decay coefficient of  $1 \times 10^{-4}$ . For the audio visual classification setting we ran 100 epochs with an initial learning rate of  $5 \times 10^{-3}$  and a mini-batch size of 32. In the audio visual continual learning task, each incremental task was also trained for 100 epochs, but with a higher learning rate of  $1 \times 10^{-2}$  and a batch size of 256 to accommodate the larger episodic memory. Audio visual question answering models were trained for 50 epochs using the same learning rate ( $1 \times 10^{-2}$ ) and a batch size of 64. On the VS100-CI benchmark, where gradient noise is significant, we replaced SGD with Adam [25] ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ) to ensure smoother convergence.

**Experimental platform.** All experiments were conducted on a linux server equipped with NVIDIA A100-40 GB GPUs and an AMD EPYC 7763 processor.

### Details of evaluation metrics

We evaluate the proposed method across three multimodal tasks using task-specific metrics.

(1) **Audio-Visual Classification.** For audio-visual classification tasks, we report the standard Top-1 accuracy,  $Acc = \frac{1}{N} \sum_{i=1}^N \mathbf{1}(\hat{y}_i = y_i)$ , where  $\mathbf{1}(\hat{y}_i = y_i)$  represents the indicator function, which equals 1 if  $\hat{y}_i = y_i$  and 0 otherwise. This metric measures the proportion of examples where the predicted label matches the ground truth.

To further assess the computational cost during training, inspired by [54], we fairly compare the efficiency

of different methods by considering both the number of epochs required to reach specified error rates and the computational complexity per epoch. Formally, the computational cost for an algorithm is defined as:  $Cost = \frac{1}{L} \sum_{l=1}^L \text{Argmin}(f(x) \leq Err_l) \times \Omega_e$ . In this formula,  $L$  represents the number of predefined error rate thresholds (set to 5 in our experiments,  $Err_l$  denotes the predefined error rate levels,  $\text{Argmin}(f(x) \leq Err_l)$  is the first epoch at which the algorithm reaches or goes below the specified error rate  $Err_l$  and  $\Omega_e$  represents the algorithmic complexity per epoch, measured in floating-point operations (FLOPs). We define the error rate  $Err_l$  using upper and lower bounds determined from the training curves of all methods under comparison. Specifically, The upper bound is set as the minimum value among the highest error rates of all compared methods. The lower bound is set to the maximum of the lowest values of the final error rates of the various methods. Within this range, we choose error rate thresholds at uniform intervals to ensure that the entire performance interval has been systematically and fairly evaluated.

(2) **Audio-Visual Continual Learning.** We track model accuracy throughout the training process using two metrics: average accuracy (AA) and average incremental accuracy (AIA). At each step  $k$  (i.e., after learning the  $k$ -th task), we compute the average accuracy  $AA_k$  across all tasks encountered so far as:  $AA_k = \frac{1}{k} \sum_{j=1}^k a_{k,j}$ , where  $a_{k,j}$  is the accuracy on the  $j$ -th task after learning task  $k$ , and  $j \leq k$ . To summarize performance across the entire learning sequence, we report the Average Incremental Accuracy, which is the mean of AA values over all  $K$  tasks:  $AIA = \frac{1}{K} \sum_{i=1}^K AA_i$ . Here, AA reflects the model’s performance after each task, while AIA captures the overall trend and stability of learning across all tasks.

To further quantify how much the model forgets previous tasks, we introduce the average forgetting rate (AFR) in the appendix Table S4. Let  $a_{k,j}$  be the test accuracy on task  $j$  after learning task  $k$ . Define the forgetting on task  $k$  as  $F_k = \frac{1}{k-1} \sum_{j=1}^{k-1} \max_{1 \leq \ell \leq k-1} (a_{\ell,j} - a_{k,j})$ , i.e., the average drop from the highest accuracy ever achieved on task  $j$  to its accuracy after the final task. For a total of  $K$  tasks, then  $AFR = \frac{1}{K-1} \sum_{k=2}^K F_k$ , where  $F_k$  excludes the first task (as forgetting can only be measured after learning at least two tasks). This metric summarizes how much the model’s performance on previously learned tasks degrades over the entire learning sequence.

(3) **Audio-Visual Question Answering.** We evaluate modal-specific question and answer accuracy, denoted by  $(A_a, A_v, A_{av})$ , which respectively evaluate performance across audio-only, visual-only, and audio-visual questions. Each accuracy is computed as  $\frac{1}{N} \sum_{i=1}^N \mathbf{1}(\widehat{\text{ans}}_i = \text{ans}_i)$ , which measures exact match between predicted and ground-truth answers over all question types within each modality.

## Code and reproducibility

The code implementation is based on Pytorch [36]. The full source code, configuration files and pre-trained checkpoints are released under an MIT licence at <https://github.com/Brain-Cog-Lab/IEMF>.

## 5. Acknowledgment

This work is supported by National Natural Science Foundation of China (NSFC) Young Scientists Fund (Grant No. 62406325).

## References

- [1] Hongjoon Ahn, Jihwan Kwak, Subin Lim, Hyeonsu Bang, Hyojun Kim, and Taesup Moon. Ss-il: Separated softmax for incremental learning. In *Proceedings of the IEEE/CVF International conference on computer vision*, pages 844–853, 2021. 6, 18
- [2] Relja Arandjelovic and Andrew Zisserman. Look, listen and learn. In *Proceedings of the IEEE international conference on computer vision*, pages 609–617, 2017. 5, 12
- [3] Marie Avillac, Sophie Deneve, Etienne Olivier, Alexandre Pouget, and Jean-René Duhamel. Reference frames for representing visual and tactile locations in parietal cortex. *Nature neuroscience*, 8(7):941–949, 2005. 2
- [4] Nick E Barraclough, Dengke Xiao, Chris I Baker, Mike W Oram, and David I Perrett. Integration of visual and auditory information by superior temporal sulcus neurons responsive to the sight of actions. *Journal of cognitive neuroscience*, 17(3):377–391, 2005. 2
- [5] Guillaume Bellec, Darjan Salaj, Anand Subramoney, Robert Legenstein, and Wolfgang Maass. Long short-term memory and learning-to-learn in networks of spiking neurons. *Advances in neural information processing systems*, 31, 2018. 17
- [6] David A Bulkin and Jennifer M Groh. Seeing sounds: visual and auditory interactions in the brain. *Current opinion in neurobiology*, 16(4):415–419, 2006. 2
- [7] Khalafalla O Bushara, Jordan Grafman, and Mark Hallett. Neural correlates of auditory–visual stimulus onset asynchrony detection. *Journal of Neuroscience*, 21(1):300–304, 2001. 2
- [8] Gemma A Calvert, Ruth Campbell, and Michael J Brammer. Evidence from functional magnetic resonance imaging of crossmodal binding in the human heteromodal cortex. *Current biology*, 10(11):649–657, 2000. 2
- [9] Huawei Cao, David G Cooper, Michael K Keutmann, Ruben C Gur, Ani Nenkova, and Ragini Verma. Crema-d: Crowd-sourced emotional multimodal actors dataset. *IEEE transactions on affective computing*, 5(4):377–390, 2014. 5, 12
- [10] Zebang Cheng, Zhi-Qi Cheng, Jun-Yan He, Kai Wang, Yuxiang Lin, Zheng Lian, Xiaojiang Peng, and Alexander Hauptmann. Emotion-llama: Multimodal emotion recognition and reasoning with instruction tuning. *Advances in Neural Information Processing Systems*, 37:110805–110853, 2024. 2
- [11] Peter Dayan and Laurence F Abbott. *Theoretical neuroscience: computational and mathematical modeling of neural systems*. MIT press, 2005. 10
- [12] Jamie Enoch, Leanne McDonald, Lee Jones, Pete R Jones, and David P Crabb. Evaluating whether sight is the most valued sense. *JAMA ophthalmology*, 137(11):1317–1320, 2019. 2
- [13] Marc O Ernst and Heinrich H Bülthoff. Merging the senses into a robust percept. *Trends in cognitive sciences*, 8(4):162–169, 2004. 1
- [14] Christopher R Fetsch, Gregory C DeAngelis, and Dora E Angelaki. Bridging the gap between theories of sensory cue integration and the physiology of multisensory neurons. *Nature Reviews Neuroscience*, 14(6):429–442, 2013. 3
- [15] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. 12
- [16] Lingyue Guo, Zeyu Gao, Jinye Qu, Suiwu Zheng, Runhao Jiang, Yanfeng Lu, and Hong Qiao. Transformer-based spiking neural networks for multimodal audio-visual classification. *IEEE Transactions on Cognitive and Developmental Systems*, 2023. 5, 12
- [17] W David Hairston, Donald A Hodges, Jonathan H Burdette, and Mark T Wallace. Auditory enhancement of visual temporal order judgment. *Neuroreport*, 17(8):791–795, 2006. 9
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 13
- [19] Shawn Hershey, Sourish Chaudhuri, Daniel PW Ellis, Jort F Gemmeke, Aren Jansen, R Channing Moore, Manoj Plakal, Devin Platt, Rif A Saurous, Bryan Seybold, et al. Cnn architectures for large-scale audio classification. In *2017 ieee international conference on acoustics, speech and signal processing (icassp)*, pages 131–135. IEEE, 2017. 13
- [20] Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. *Advances in Neural Information Processing Systems*, 35:28708–28720, 2022. 13
- [21] Chengpeng Jiang, Jiaqi Liu, Yao Ni, Shangda Qu, Lu Liu, Yue Li, Lu Yang, and Wentao Xu. Mammalian-brain-inspired neuromorphic motion-cognition nerve achieves cross-modal perceptual enhancement. *Nature Communications*, 14(1):1344, 2023. 2
- [22] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 12
- [23] Nitish Shirish Keskar, Jorge Nocedal, Ping Tak Peter Tang, Dheevatsa Mudigere, and Mikhail Smelyanskiy. On large-batch training for deep learning: Generalization gap and sharp minima. In *5th International Conference on Learning Representations, ICLR 2017*, 2017. 12

- [24] Minsu Kim, Joanna Hong, Se Jin Park, and Yong Man Ro. Cromm-vsr: Cross-modal memory augmented visual speech recognition. *IEEE Transactions on Multimedia*, 24:4342–4355, 2021. 2
- [25] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. 13
- [26] Guangyao Li, Yake Wei, Yapeng Tian, Chenliang Xu, Ji-Rong Wen, and Di Hu. Learning to answer questions in dynamic audio-visual scenarios. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19108–19118, 2022. 7, 13, 18
- [27] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017. 6, 18
- [28] Fengmao Lv, Xiang Chen, Yanyong Huang, Lixin Duan, and Guosheng Lin. Progressive modality reinforcement for human multimodal emotion recognition from unaligned multimodal sequences. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2554–2562, 2021. 2
- [29] Emiliano Macaluso and Jon Driver. Multisensory spatial interactions: a window onto functional integration in the human brain. *Trends in neurosciences*, 28(5):264–271, 2005. 2
- [30] Emiliano Macaluso, Nathalie George, Ray Dolan, Charles Spence, and Jon Driver. Spatial and temporal factors during processing of audiovisual speech: a pet study. *Neuroimage*, 21(2):725–732, 2004. 2
- [31] Youssef Mroueh, Etienne Marcheret, and Vaibhava Goel. Deep multimodal learning for audio-visual speech recognition. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2130–2134. IEEE, 2015. 2
- [32] Arsha Nagrani, Shan Yang, Anurag Arnab, Aren Jansen, Cordelia Schmid, and Chen Sun. Attention bottlenecks for multimodal fusion. *Advances in Neural Information Processing Systems*, 34:14200–14213, 2021. 2
- [33] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 807–814, 2010. 10
- [34] Toemme Noesselt, Jochem W Rieger, Mircea Ariel Schoenfeld, Martin Kanowski, Hermann Hinrichs, Hans-Jochen Heinze, and Jon Driver. Audiovisual temporal correspondence modulates human multisensory superior temporal sulcus plus primary sensory cortices. *Journal of Neuroscience*, 27(42):11431–11441, 2007. 2
- [35] Uta Noppeney. Perceptual inference, learning, and attention in a multisensory world. *Annual review of neuroscience*, 44:449–473, 2021. 1
- [36] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Razion, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. 14
- [37] Xiaokang Peng, Yake Wei, Andong Deng, Dong Wang, and Di Hu. Balanced multimodal learning via on-the-fly gradient modulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8238–8247, 2022. 2, 5
- [38] Weiguo Pian, Shentong Mo, Yunhui Guo, and Yapeng Tian. Audio-visual class-incremental learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7799–7811, 2023. 6, 12, 18
- [39] Christina Regenbogen, Janina Seubert, Emilia Johansson, Andreas Finkelmeyer, Patrik Andersson, and Johan N Lundström. The intraparietal sulcus governs multisensory integration of audiovisual information based on task difficulty. *Human brain mapping*, 39(3):1313–1326, 2018. 2, 9
- [40] Herbert Robbins and Sutton Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951. 13
- [41] Muhtasim Ul Karim Sadaf, Najam U Sakib, Andrew Pannone, Harikrishnan Ravichandran, and Saptarshi Das. A bio-inspired visuotactile neuron for multisensory integration. *Nature Communications*, 14(1):5729, 2023. 2
- [42] Daniel Senkowski, Dave Saint-Amour, Thomas Gruber, and John J Foxe. Look who’s talking: the deployment of visuo-spatial attention during multisensory speech processing under noisy environmental conditions. *Neuroimage*, 43(2):379–387, 2008. 2
- [43] Daniel A Slutsky and Gregg H Recanzone. Temporal and spatial dependency of the ventriloquism effect. *Neuroreport*, 12(1):7–10, 2001. 9
- [44] Barry E Stein and M Alex Meredith. *The merging of the senses*. MIT press, 1993. 9
- [45] Barry E Stein and Terrence R Stanford. Multisensory integration: current issues from the perspective of the single neuron. *Nature reviews neuroscience*, 9(4):255–266, 2008. 3
- [46] Barry E Stein, Terrence R Stanford, and Benjamin A Rowland. Development of multisensory integration from the perspective of the individual neuron. *Nature Reviews Neuroscience*, 15(8):520–535, 2014. 9
- [47] Gregor Rafael Szycik, Peggy Tausche, and Thomas F Münte. A novel approach to study audiovisual integration in speech perception: localizer fmri and sparse sampling. *Brain Research*, 1220:142–149, 2008. 2
- [48] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. *Advances in neural information processing systems*, 35:10078–10093, 2022. 13
- [49] Nienke Van Atteveldt, Elia Formisano, Rainer Goebel, and Leo Blomert. Integration of letters and speech sounds in the human brain. *Neuron*, 43(2):271–282, 2004. 2

- [50] Yang Yang, Fengqiang Wan, Qing-Yuan Jiang, and Yi Xu. Facilitating multimodal classification via dynamically learning modality gap. *Advances in Neural Information Processing Systems*, 37:62108–62122, 2024. 5
- [51] Yiqun Yao and Rada Mihalcea. Modality-specific learning rates for effective multimodal additive late-fusion. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1824–1834, 2022. 5
- [52] Jeong Hun Yeo, Minsu Kim, Jeongsoo Choi, Dae Hoe Kim, and Yong Man Ro. Akvsr: Audio knowledge empowered visual speech recognition by compressing audio knowledge of a pretrained model. *IEEE Transactions on Multimedia*, 2024. 2
- [53] Fangwen Yu, Yujie Wu, Songchen Ma, Mingkun Xu, Hongyi Li, Huanyu Qu, Chenhang Song, Taoyi Wang, Rong Zhao, and Luping Shi. Brain-inspired multimodal hybrid neural network for robot place recognition. *Science Robotics*, 8(78):eabm6996, 2023. 2
- [54] Tielin Zhang, Xiang Cheng, Shuncheng Jia, Mu-ming Poo, Yi Zeng, and Bo Xu. Self-backpropagation of synaptic modifications elevates the efficiency of spiking and artificial neural networks. *Science advances*, 7(43):eabh0146, 2021. 13

## A. appendix

Category	Parameters	Values
Audio Visual Classification	Network backbone	ResNet-18
	Optimizer	SGD
	Weight decay	$1 \times 10^{-4}$
	Initial learning rate	$5 \times 10^{-3}$
	Number of training epochs	100
	Batch size	32
Audio Visual Continual Learning	Learning rate	$1 \times 10^{-2}$
	Batch size	256
	Number of training epochs	100
	Optimizer (VS100-CI)	Adam ( $\beta_1 = 0.9, \beta_2 = 0.999$ )
Audio Visual Question and Answering	Learning rate	$1 \times 10^{-2}$
	Batch size	64
	Number of training epochs	50
LIF Neuron (SNN)	Resting potential $V_{\text{rest}}$	0
	Firing threshold $V_{\text{th}}$	0.5
	Membrane time constant $\tau_m$	2.0
	Surrogate gradient function	Piecewise linear [5]
	Conductivity $g_L, g_E, g_I$	1 S, 1 S, 1 S
	Reversal potential $V_E, V_I$	0
	Discrete time step $t$	4
IEMF	Inverse gain coefficient $\gamma$	1.0
	Gating function $\kappa(\cdot)$	$\tanh(\cdot)$

Table S1. Comprehensive experimental parameters used for implementation across three multimodal tasks.

Methods	CREMA-D				Kinetics-Sounds				UrbanSound8K-AV			
	Normal	MSLR	OGM_GE	LFM	Normal	MSLR	OGM_GE	LFM	Normal	MSLR	OGM_GE	LFM
Vanilla	62.63	64.11	68.68	<b>64.11</b>	51.58	51.89	57.63	55.28	<b>97.90</b>	97.79	97.60	98.05
w/ IEMF	<b>63.44</b>	<b>65.59</b>	<b>71.10</b>	63.98	<b>56.17</b>	<b>55.86</b>	<b>64.61</b>	<b>63.15</b>	97.86	<b>97.98</b>	<b>99.24</b>	<b>98.63</b>

Table S2. Comparison of the proposed method (w/ IEMF) and the vanilla baseline on ANN across three multimodal datasets—CREMA-D, Kinetics-Sounds, and UrbanSound8K-AV—under four fusion methods (Normal, MSLR, OGM\_GE, and LFM). Bold values indicate the highest accuracy achieved in each configuration.

Methods	CREMA-D				Kinetics-Sounds				UrbanSound8K-AV			
	Normal	MSLR	OGM_GE	LFM	Normal	MSLR	OGM_GE	LFM	Normal	MSLR	OGM_GE	LFM
Vanilla	63.04	63.70	69.49	63.58	53.12	53.19	57.28	54.63	<b>98.02</b>	97.98	97.67	98.13
w/ IEMF	<b>64.65</b>	<b>64.38</b>	<b>69.35</b>	<b>64.78</b>	<b>54.97</b>	<b>56.63</b>	<b>62.88</b>	<b>63.53</b>	97.79	<b>98.05</b>	<b>99.35</b>	<b>98.63</b>

Table S3. Same as Table S2, with the difference that evaluation is conducted on SNN.

Method	AVE-CI		K-S-CI		VS100-CI	
	Mean Accuracy ↑	Avg Forgetting ↓	Mean Accuracy ↑	Avg Forgetting ↓	Mean Accuracy ↑	Avg Forgetting ↓
LwF [27]	54.06	26.77	59.89	<b>15.26</b>	59.06	<b>17.91</b>
LwF w/ IEMF	<b>56.00</b>	<b>24.37</b>	<b>62.79</b>	17.25	<b>60.72</b>	18.25
SSIL [1]	56.19	7.61	63.31	<b>4.66</b>	68.21	<b>8.53</b>
SSIL w/ IEMF	<b>58.12</b>	<b>6.25</b>	<b>65.18</b>	5.45	<b>70.01</b>	9.30
AV-CIL [38]	61.86	24.31	70.63	11.03	70.15	<b>8.98</b>
AV-CIL w/ IEMF	<b>61.94</b>	<b>20.87</b>	<b>72.49</b>	<b>10.44</b>	<b>70.76</b>	9.49

Table S4. Performance comparison on three audio-visual continual learning benchmarks (AVE-CI, K-S-CI, and VS100-CI) in terms of Mean Accuracy ( $\uparrow$ ) and Average Forgetting ( $\downarrow$ ). We evaluate three baseline methods (LwF, SSIL, and AV-CIL) and their variants augmented with the proposed IEMF module. Bold values indicate the best performance in each metric.

Method	Audio Question (%)			Visual Question (%)			Audio-Visual Question (%)					Overall Avg. (%)	
	Counting	Comparative	Avg	Counting	Location	Avg	Existential	Location	Counting	Comparative	Temporal	Avg	
Baseline	77.20	62.06	71.60	74.15	<b>76.79</b>	75.48	81.71	67.43	<b>62.57</b>	61.61	62.99	67.34	70.24
Baseline w/ IEMF	<b>77.40</b>	<b>63.89</b>	<b>72.40</b>	<b>75.23</b>	<b>76.79</b>	<b>76.02</b>	<b>82.11</b>	<b>69.07</b>	59.55	<b>62.87</b>	<b>64.93</b>	<b>67.87</b>	<b>70.82</b>
ST-AVQA [26]	77.59	62.23	71.90	73.89	75.57	74.74	<b>82.81</b>	68.45	<b>63.00</b>	60.45	62.86	67.61	70.26
ST-AVQA w/ IEMF	<b>79.84</b>	<b>65.39</b>	<b>74.49</b>	<b>74.65</b>	<b>76.63</b>	<b>75.65</b>	82.11	<b>69.47</b>	62.68	<b>62.51</b>	<b>64.20</b>	<b>68.33</b>	<b>71.36</b>

Table S5. Comparison results with different AVQA methods on the MUSIC-AVQA dataset, where different types of questions (audio-only, visual-only, and audio-visual) are evaluated.