# Music's Multimodal Complexity in AVQA: Why We Need More than General Multimodal LLMs

Wenhao You<sup>1</sup>, Xingjian Diao<sup>2</sup>, Chunhui Zhang<sup>2</sup>, Keyi Kong<sup>3</sup>, Weiyi Wu<sup>2</sup>, Zhongyu Ouyang<sup>2</sup>, Chiyu Ma<sup>2</sup>, Tingxuan Wu<sup>4</sup>, Noah Wei<sup>5</sup>, Zong Ke<sup>6</sup>, Ming Cheng<sup>2</sup>, Soroush Vosoughi<sup>2</sup>, Jiang Gui<sup>2</sup> <sup>1</sup>University of Waterloo, <sup>2</sup>Dartmouth College, <sup>3</sup>Shandong University, <sup>4</sup>The London School of Economics and Political Science, <sup>5</sup>University of Hong Kong, <sup>6</sup>National University of Singapore

### Abstract

While recent Multimodal Large Language Models exhibit impressive capabilities for general multimodal tasks, specialized domains like music necessitate tailored approaches. Music Audio-Visual Question Answering (Music AVQA) particularly underscores this, presenting unique challenges with its continuous, densely layered audio-visual content, intricate temporal dynamics, and the critical need for domain-specific knowledge. Through a systematic analysis of Music AVQA datasets and methods, this position paper identifies that specialized input processing, architectures incorporating dedicated spatial-temporal designs, and music-specific modeling strategies are critical for success in this **domain.** Our study provides valuable insights for researchers by highlighting effective design patterns empirically linked to strong performance, proposing concrete future directions for incorporating musical priors, and aiming to establish a robust foundation for advancing multimodal musical understanding. This work is intended to inspire broader attention and further research, supported by a continuously updated anonymous GitHub repository of relevant papers: https://github.com/xid32/Survey4MusicAVQA.

### **1** Introduction

"Music is a moral law. It gives a soul to the Universe, wings to the mind, flight to the imagination, a charm to sadness, gaiety and life to everything."

- Plato (c. 427-347 BCE, Ancient Greece)

Multimodal Large Language Models (MLLMs) have demonstrated impressive efficacy across a wide range of tasks, modelling various modalities such as text, image, and audio [1, 2, 3, 4, 5, 6]. However, this success brings forth a critical consideration: the tension between the broad applicability of general multimodal approaches and the requirements of specialized domains [7, 8, 9]. This leads to a central question: Are general-purpose MLLMs, despite their advancements, truly sufficient for all multimodal tasks, especially those demanding deep, domain-specific understanding?

Music Performance Audio-Visual Question Answering (Music AVQA) emerges as a particularly challenging multimodal domain that compellingly illustrates this tension [10, 11, 12, 13]. Music, in its rich complexity, often requires specialized treatment beyond the capabilities of generic models [14, 15, 16, 17, 18, 19, 20, 21, 22]. Unlike common scenarios with sparse and discrete audio signals, music performances exhibit a continuous and tightly interwoven blend of audio and visual signals, offering a uniquely rich context for fine-grained audio-visual scene understanding and temporal reasoning [10, 11, 14, 23, 24, 25]. For instance, tasks such as discerning the loudest instrument



Figure 1: Contrast between (i) conventional QA and (ii) Music-AVQA with dense audio. Panel (i) shows an isolated sound (barking) and synchronized action, which are relatively easy to detect. Panel (ii) exemplifies music's complexity, featuring overlapping instruments and rhythmic patterns. Such dense and continuous audio-visual signals demand fine-grained temporal and spatial reasoning through cross-modal comparisons and they are more challenging than conventional multimodal QA.

amidst an ensemble or comparing rhythmic complexity between two spatially distinct performers demand a level of granularity that general-purpose MLLMs may not inherently possess [9, 14, 24]. Thus, Music AVQA serves as an ideal lens through which to examine the limitations of general multimodal approaches and to advocate for the necessity of domain-specific adaptations [10, 14].

The unique challenges of Music AVQA, illustrated in Figure 1, stem from the need to reason over continuous, temporally evolving, and densely layered audio-visual signals. Specific complexities include: *First*, musical pieces often contain dense and layered audio information. Multiple overlapping instrumental sources are common, which necessitates fine-grained processing to disentangle and interpret complex auditory scenes. *Second*, effective understanding depends on precise temporal alignment. It is crucial to accurately associate visual cues—such as a musician's actions—with their corresponding auditory outputs. This alignment must occur across multiple timescales and often involves intricate temporal dynamics. *Third*, the domain frequently requires specialized knowledge. This includes instrument recognition, familiarity with musical theory (such as rhythm and harmony), and an understanding of performance conventions, whether these are explicit or implicit. *Finally*, Music AVQA questions often involve complex spatial-temporal relationships. For example, one may need to track dynamic intensity across simultaneous sources ("Which instrument produces the loudest sound?") or reason about spatial and temporal rhythmic patterns ("Is the cello on the right more rhythmic than the cello on the left?"). Collectively, these factors underscore the unique and demanding nature of reasoning in Music AVQA.

This study argues that Music AVQA is a fundamentally distinct multimodal reasoning task, for which specialized multimodal designs are essential and empirically linked to strong model performance. As the first comprehensive survey in this area, we specifically analyze how tailored designs—particularly in input processing and spatial-temporal architecture—enable more effective music understanding compared to generic multimodal systems. Furthermore, we outline how such specialized approaches, by incorporating deeper musical priors, can further advance the field.

## 2 Background

Why Music AVQA is more challenging than normal multimodal understanding? Music AVQA presents several distinctive challenges: ① Dense Signal Interpretation: Unlike sparse audio events in conventional AVQA, music performances feature continuous, overlapping instrumental sources that require sophisticated separation and attribution; ② Hierarchical Temporal Reasoning: Musical information unfolds across multiple time scales (beats, phrases, sections), demanding models capable of reasoning across these hierarchical structures; ③ Cross-Modal Correspondence: Establishing reliable associations between visual instrumental actions and their acoustic outputs is complicated by temporal misalignments between physical gestures and the resulting sounds; ④ Domain-Specific Knowledge: Effective reasoning often depends on implicit musical knowledge, such as instrumental techniques, ensemble conventions, and acoustic properties; ⑤ Abstract Attribute Quantification: Questions involving subjective qualities such as "rhythmic", "melodic," or "harmonious" require computational strategies to map linguistic descriptors onto measurable signal properties; ⑥ Data Scarcity: The specialized nature of musical performances results in smaller and less diverse datasets compared to general AVQA tasks, limiting the generalization capabilities of trained models.

What are common music performance scene types? ① Solo Performance – A single musician showcasing technical skills and artistic expression on one instrument. ② Ensemble of the Same Instrument – Multiple musicians playing identical or related instruments, creating unified harmonies and textures. ③ Ensemble of Different Instruments – Musicians performing with a variety of instruments, producing diverse tonal colors and complex musical interactions. ④ Culture-Specific Ensemble – Traditional instrumental groups that embody the musical heritage and regional styles of specific cultures. See Appendix Section B for examples.

What are common question types in Music AVQA? ① Existential Questions: Determine whether a sound corresponds to a visible object in the scene (e.g., "Is this sound from the instrument in the video?"). ② Counting Questions: Quantify audio-visual elements that require cross-modal integration (e.g., "How many instruments are sounding in the video?"). ③ Location Questions: Identify the spatial position of sound sources within the visual scene (e.g., "Where is the first sounding instrument?"). ④ Comparative Questions: Compare properties across different audio-visual elements (e.g., "Is the instrument on the left louder than the one on the right?"). ⑤ Temporal Questions: Reason about the timing and sequential relationships between auditory and visual events (e.g., "Which instrument produces sound before the piano?"). See Appendix Section C for examples.

## **3** Evolution of MUSIC-AVQA Datasets

The development of Music AVQA research has been driven by progressively refined datasets addressing specific limitations. As summarized in Table 6 in Appendix Section E, this evolution began with the ① MUSIC-AVQA dataset [10], the first large-scale benchmark designed specifically for AVQA in musical contexts, comprising 9,288 performance videos and 45,867 question-answer pairs across diverse reasoning tasks. Subsequent research reveal challenges related to data bias and imbalanced answer distributions, prompting the creation of ② MUSIC-AVQA v2.0 [11], which expands to 10,518 videos and approximately 54,000 question-answer pairs. This version balance 15 biased templates by ensuring no dominant answers exceed 60% for binary questions or 50% for multi-class questions, particularly enhancing representation in various question categories. Building on these foundations, ③ MUSIC-AVQA-R [23] introduces robustness evaluation through question rephrasing, expanding the test set from 9,129 to 211,572 questions. With a vocabulary five times larger than the original dataset, MUSIC-AVQA-R distinguishes between head (common) and tail (rare) samples, enabling assessment of model performance in both in-distribution and out-of-distribution scenarios. This progressive refinement of datasets has laid a solid foundation for advancing multimodal understanding and robust evaluation in music performance environments.

### 4 Categorization of Music AVQA Methods Based on Architecture

Music AVQA methods exhibit diverse architectural designs, particularly in how they encode and integrate textual, visual, and auditory modalities. To better organize existing approaches by their core modeling strategies, we categorize them into three groups—Transformer-based, CNN-based, and Hybrid models—as summarized in Table 1. This categorization highlights how different models are structured to handle the continuous and densely layered nature of musical performances.

**Transformer-based models.** Transformer-based models are characterized by the extensive use of self-attention mechanisms, which benefit in particular from their ability to handle long-range temporal dependencies and fine-grained cross-modal alignment. Methods such as Amuse utilize transformers across all modalities, combining a Swin Transformer for visual processing with an HTS-AT transformer for audio encoding, and employing cross-modal adapters to facilitate early and frequent fusion of multimodal information. Similarly, LAST-Att integrates a Swin-V2 Transformer for vision and an Audio Spectrogram Transformer (AST) for audio, emphasizing fine-grained spatial-temporal alignment through pixel-level cross-modal attention. Other methods such as LAVisH and LSTTA, adopt lightweight transformer adapters to inject multimodal cues into frozen transformer backbones, enabling efficient cross-modal reasoning while leveraging strong pre-trained representations.

**CNN-based models.** CNN-based methods typically utilize convolutional backbones such as ResNet or VGGish to encode modality-specific information into global or regional features, often relying on simpler late-stage fusion strategies. The AVST method exemplifies this approach, combining ResNet-18 visual embeddings and VGGish audio features through spatial attention modules to

Метнор	Text Encoder	Visual Encoder	Audio Encoder	S-T
Amuse [14]	Transformer [26]	Swin-Transformer-v2 [27]	HTS-AT [28]	1
AUDIO FLAMINGO [29]	OPT-IML-MAX-1.3B [30]	-	ClapCap [31]	1
AVMOE [32]	-	Swin-Transformer-v2 [27]	HTS-AT [28]	×
AVSD [33]	LSTM	LSTM	LSTM	×
AVSIAM [34]	-	ViT [35]	ViT [35]	×
AVST [10]	LSTM	ResNet-18 [36]	VGGish [37]	1
CAT [13]	LLaMA2-7B [38]	ImageBind [1]	ImageBind [1]	×
CHATBRIDGE [39]	Vicuna-13B [40]	ViT-G [41]	BEATs [42]	×
CIGN [43]	-	ResNet-18 [36]	ResNet-18 [36]	1
COCA [44]	Word Embedding	ResNet-18 [36]	VGGish [37]	×
CONVLSTM [45]	LSTM	-	Conv	×
CROSSMAE [46]	-	MAE [47]	AudioMAE [48]	×
DCL [49]	DeBERTa-V3-Large [50]	ViT [35]	AST [51]	1
DG-SCT [52]	-	ViT [35]	HTS-AT [28]	1
EEMC [53]	RoBERTa [54]	ViT [35]	VGGish [37]	1
FCNLSTM [45]	LSTM	-	Conv	×
GPT-40 [4]	Transformer	CLIP-ViT	Transformer	×
GRU [55]	LSTM	VGGNet [56]		×
HCRN [57]	BiLSTM	ResNet-18 [36]	-	×
LAST-ATT [11]	LSTM	Swin-Transformer-v2 [27]	Audio-Spectrogram-Transformer	1
LAVISH [58]		ViT [35]	ViT [35]	1
LAVIT [59]	Transformer [26]	Transformer [26]	Transformer [26]	1
LSTTA [60]	CLIP [35]	CLIP [35]	w2v-Conformer [61]	1
MAVEN [62]	Mixtral	InternViT-300M-448px [63]	Transformer	×
MCAN [64]	GloVe [65]+LSTM	Faster R-CNN [66]		×
MCCD [23]	-	-	-	1
MEERKAT [67]	LLaMA2-7B [38]	CLIP-ViT	CLAP [68]	1
OGM [69]	-	ResNet-18 [36]	ResNet-18 [36]	×
ONELLM [2]	LLaMA2-7B [38]	CLIP-ViT	Unified Multimodal Encoder	×
OPM [69]	-	ResNet-18 [36]	ResNet-18 [36]	×
PSAC [70]	Word Embedding	CNN	-	×
PSTP-NET [25]	CLIP [35]	CLIP [35]	VGGish [37]	1
OAP [5]	DeBERTa-V2-XLarge	CLIP [35]	CLAP [68]	×
OWEN2.5-VL [3]	MRoPE [3]	ViT [35]		×
REFATOMNET [71]	BERT	ViT [35]	-	1
VALOR [72]	BERT	CLIP [35]	AST [51]	×
VAST [73]	BERT [74]	ViT [75]	BEATs [42]	×
VIDEOLLAMA-2 [76]	Transformer	CLIP [35]	BEATs [42]	1
VITA [77]	Mixtral [78]	InternViT-300M-448px [63]	CNN	×

Table 1: Architectural summary of representative Music AVQA methods. Each method lists the text, visual, and audio encoders used, along with an indication of whether explicit spatial-temporal (S-T) modeling is incorporated. Detailed descriptions of each method are provided in Appendix F and G.

explicitly localize sound sources within visual frames. PSTP-Net extends this design by introducing a progressive refinement strategy that sequentially filters temporal segments and spatial regions, systematically narrowing down question-relevant audio-visual content prior to fusion. Although CNN-based models are computationally efficient and straightforward, their reliance on late fusion may pose challenges to capturing the complex temporal dynamics characteristic of musical performances.

**Hybrid models.** Hybrid models combine CNNs, transformers, and large language models (LLMs) to enable unified multimodal reasoning. They typically employ pre-trained encoders from both CNN and transformer families, integrated through sophisticated cross-modal fusion mechanisms. Representative examples include ChatBridge, CAT, OneLLM, and Meerkat. ChatBridge utilizes a perceiver-based multimodal transformer to merge modalities via language-aligned latent representations, followed by a frozen LLM for reasoning. CAT introduces modality-specific clue aggregation modules on top of ImageBind encodings, enabling precise question-driven multimodal grounding before passing information to a generative LLaMA2 LLM. OneLLM further generalizes multimodal integration by introducing a universal projection mechanism that allows a single LLM to interpret diverse modality embeddings seamlessly. In contrast, Meerkat emphasizes fine-grained cross-modal alignment through an audio-visual optimal transport module that explicitly matches audio segments to corresponding visual regions, achieving strong performance on tasks requiring precise localization of sound sources, underscoring the benefit of precise local grounding for complex audio-visual interactions in musical contexts.

## 5 A Call on Specialized Multimodal Input Processing for Music AVQA

While input preparation is often treated as a fixed pipeline in general AVQA, music performance settings introduce unique challenges that make input fidelity, segmentation, and representation design

especially consequential. Musical scenes are densely layered, temporally continuous, and rich in expressive detail, requiring greater care in how audio, visual, and textual inputs are captured and structured. In what follows, we examine how Music AVQA tasks motivate specialized input processing across three key fronts: maintaining high-resolution and synchronized multimodal signals, adapting tokenization to the structure of musical content, and managing the scale and diversity of music-specific data representations.

**Continuous, high-fidelity, and tightly aligned inputs are foundational.** Compared to eventcentric AVQA tasks that typically involve short, discrete sound events and lower-resolution recordings, Music AVQA deals with continuous, polyphonic streams spanning multiple spatial and temporal scales. Audio is commonly sampled at high rates (44.1 kHz or above) and often preserved in lossless formats to retain subtle timbral and articulatory detail [79]. Visual inputs similarly tend to require higher resolution (1080p or above) and frame rates (30–60 fps) to capture nuanced performer motions such as bowing or fingering [80, 81]. Even modest temporal offsets—around 100–200ms—can affect the perceived correspondence between gesture and sound. To improve synchronization and cue isolation, some recent models adopt preprocessing strategies like beat-based segmentation [82] and harmonic-percussive separation [83], which can help surface rhythmically or acoustically meaningful content for downstream reasoning.

**Tokenization strategies benefit from musical adaptation.** Tokenization plays a central role in structuring inputs for multimodal reasoning, and recent Music AVQA models often tailor their strategies to preserve musical structure. For audio, models such as AMUSE [14], DG-SCT [52], and PSTP-NET [25] transform waveforms into Mel-spectrograms, which are then segmented via patch-based encoders like AST [51] and HTS-AT [28] or CNNs such as VGGish [37] and ResNet-18 [36]. AUDIO FLAMINGO [29], for instance, uses overlapping 7-second windows in CLAPCAP [31] to embed long-range audio context. Visual streams are frequently tokenized using ViT [35] or Swin-based [27] patch embeddings (e.g., in AVSIAM [34] and LAVISH [58]), while earlier models like AVST [10] use frame-level CNN features. Text tokenization is typically handled by subword models aligned with large language models (e.g., LLAMA2 [38], ROBERTA [54]), as seen in CHATBRIDGE [39] and ONELLM [2]. These tokenization schemes help preserve temporal granularity and modality alignment, which may be important for interpreting overlapping instruments, rhythmic changes, and localized visual cues.

**Musical content introduces distinct data and representational considerations.** Music AVQA tasks often involve long-form performances with overlapping sources and evolving musical dynamics, which can create challenges for segmentation, annotation, and generalization. Unlike typical AVQA datasets centered on short clips and isolated actions, music-focused benchmarks (e.g., MUSIC-AVQA [10]) include multi-instrument performances spanning several minutes. These conditions place greater demands on dataset diversity to avoid overfitting to genre-specific patterns or ensemble configurations. To broaden coverage, some models are trained on data drawn from live performances, studio recordings, and synthetic renderings. However, the absence of symbolic structure can limit the model's access to mid-level grounding. In this context, musically informed preprocessing (e.g., onset alignment, rhythmic segmentation, graph representation learning [84]) may support more interpretable and temporally aligned input representations.

## 6 A Call on Specialized Spatial-Temporal Designs for Music AVQA

We systematically analyze the models listed in Table 1 to identify architectural factors associated with strong Music AVQA performance across diverse multimodal designs. Each model is annotated based on whether it incorporates **spatial-temporal design**, defined as architectural components explicitly aimed at localizing audio-visual content in space and time—such as temporal segment selection, spatial attention, or cross-modal alignment modules. This categorization enables us to assess whether high-performing models exhibit structural traits aligned with the temporally continuous and spatially layered nature of musical performances.

To assess the empirical impact of spatial-temporal design, we evaluate Music AVQA models across representative question types grouped by modality—audio, visual, and audio-visual—as shown in Figure 2. Each subplot compares model accuracy on a specific QA type, with bars color-coded to indicate

whether spatial-temporal design is applied for the relevant modality. This setup allows precise attribution of performance differences to design choices. To capture broader trends, Figure 3 summarizes average accuracy across all 13 QA categories using radar plots on two benchmarks: Music-AVQA and Music-AVQA-R. These visualizations reveal that models with spatial-temporal design consistently outperform their counterparts, particularly in tasks involving fine-grained localization or temporal sequencing. The full quantitative results supporting these figures are reported in Appendix A, Tables 2, 3, and 4. This experimental design enables systematic assessment of spatial-temporal design as a key architectural driver of multimodal reasoning in musical environments.



(g) Audio-Visual Comparative QA. (h) Audio-Visual Temporal QA.

(i) Audio-Visual QA Average.

Figure 2: Accuracy comparison of Music AVQA models across representative question types, grouped by modality: (a–b) Audio, (c–e) Visual, and (f–i) Audio-Visual. Each bar corresponds to a model and is color-coded based on whether it incorporates **spatial-temporal design** for the relevant task type: bars in green, purple, and orange represent models that apply spatial-temporal modeling to Audio-related, Visual-related, and Audio-Visual-related question answering, respectively; bars in blue represent models without spatial-temporal design. Across most categories, models with spatial-temporal components tend to perform more accurately, particularly on tasks requiring temporal reasoning or spatial localization. These patterns suggest that incorporating spatial-temporal design supports more effective reasoning in musically structured multimodal environments.

Spatial-temporal design enhances audio QA by supporting fine-grained tracking of overlapping sources and temporally evolving acoustic cues. Audio-related questions in Music AVQA—such as instrument counting or loudness comparison—require models to distinguish simultaneous sound sources, localize temporal onsets, and resolve dynamic variations across time. As shown in Figures 2(a) and 2(b), models with spatial-temporal design consistently outperform others. LAST-ATT [11] achieves the highest audio counting accuracy at 85.71%, benefiting from repeated cross-attention between question-guided Swin-Transformer features and spectrogram patches from an Audio



(a) Methods on Music-AVQA [10].

(b) Methods on Music-AVQA-R [23].

Figure 3: Radar plots showing the per-type average accuracy of model groups with and without **spatial-temporal design** across 13 QA categories on (a) Music-AVQA [10] and (b) Music-AVQA-R [23]. Each axis corresponds to a QA type spanning audio, visual, and audio-visual reasoning, including the overall average (Total-Average). The filled green polygon in Figure 3(a) and purple polygon in Figure 3(b) represent the mean accuracy across QA types for models with spatial-temporal design, while the blue polygon represents the average performance of models without such design. Models with spatial-temporal design consistently achieve higher accuracy across all modality groups. These advantages persist under distribution shift in the robustness-focused Music-AVQA-R dataset.

Spectrogram Transformer, which helps the model focus on musically salient moments. AMUSE [14], with 83.58% average audio QA accuracy, aligns audio-video streams using beat-synchronous features and temporally-adaptive fusion modules, allowing it to isolate relevant auditory content even under polyphonic conditions. DG-SCT [52] further introduces bidirectional attention layers across temporal, spatial, and channel dimensions, dynamically adjusting audio-visual focus based on the question's semantics. By contrast, models lacking spatial-temporal structure—such as MCAN (67.47%) and CONVLSTM (66.73%)—often rely on global feature pooling or frame-agnostic fusion, making them vulnerable to overlap, misalignment, and temporal drift. Notably, spatial-temporal designs adopt recurring architectural motifs: temporal segment selection (PSTP-NET [25], AVST [10]), audio-guided visual attention (DG-SCT, LSTTA [60]), and fine-grained cross-modal alignment (MEERKAT [67]). These mechanisms are well-suited for modeling music's complex structure, where overlapping instruments and evolving rhythms require localized reasoning in both time and space. The strong performance of spatial-temporal models across audio QA tasks confirms their value in resolving multi-instrument scenarios and detecting temporally grounded acoustic attributes.

Spatial-temporal design improves visual QA by enhancing spatial disambiguation and capturing motion cues over time. Visual-related questions in Music AVOA—such as counting instruments or identifying positions—often involve tracking multiple performers, detecting visual cues of articulation (e.g., bowing, striking), and resolving spatial relationships within densely packed frames. As shown in Figures 2(c)-2(e), models with spatial-temporal components generally achieve stronger accuracy. For example, LSTTA [60] (82.03% visual QA average) combines short-term semantic interaction and long-term semantic filtering modules to capture both local gestures and global scene dynamics, enabling precise reasoning about when and where instruments are engaged. DG-SCT [52] (82.08%) uses cross-modal temporal attention guided by audio prompts to enhance visual token selection, focusing on visually active regions corresponding to sounding instruments. PSTP-NET [25] (77.26%) implements a region refinement module that explicitly filters visual patches within question-relevant segments, improving spatial disambiguation. While spatial-temporal modeling is effective, some models without it still perform competitively-most notably CAT [13] (86.10%), which leverages large-scale pretrained vision encoders (ImageBind) and LLaMA2 to infer structure implicitly. However, such models may rely heavily on correlation learned from pretraining, rather than explicit reasoning about visual dynamics. Spatial-temporal models, by contrast, explicitly model the temporal unfolding of gestures and the spatial focus of performer activity—important properties in musical scenes where instrument positions are static but their activation varies over time. These architectural patterns help stabilize attention and reduce confusion when multiple instruments are

visually present but only some are active, contributing to more consistent visual QA performance across counting and localization tasks.

Spatial-temporal design is critical for audio-visual OA, where accurate reasoning requires precise temporal and spatial alignment between modalities. Among all Music AVQA categories, audio-visual questions impose the strongest demand on cross-modal synchronization, requiring the model to associate specific acoustic events with their visual sources over time. As shown in Figures 2(f)-2(i) and Table 2, models with spatial-temporal components consistently achieve higher accuracy across AV-Existential, AV-Counting, AV-Location, AV-Comparative, and AV-Temporal types. AMUSE [14] reaches 82.43% on overall AV questions by leveraging segment-level alignment between synchronized beat-level audio and video inputs and applying cross-modal adapters at each step. PSTP-NET [25] adopts a progressive three-stage pipeline: temporal segment selection, spatial region refinement, and audio-guided attention, resulting in 72.57% AV average. MEERKAT [67] further enhances local alignment by explicitly modeling cross-modal transport between audio patches and visual regions, and enforces bounding box constraints for grounding, yielding strong performance on AV-Comparative and AV-Location. In contrast, models without spatial-temporal design-such as MCAN (57.80%), GPT-40 (50.08%), and QWEN2.5-VL (47.75%)-struggle to resolve fine-grained multimodal relationships. While CAT [13] achieves 83.20% AV average through large-scale pretrained encoders, its performance drops on AV-Temporal and AV-Location tasks that require precise temporal ordering or spatial binding. These results support that spatial-temporal designs-especially those involving temporally segmented reasoning, audio-guided spatial focus, and per-frame fusion-enable the model to track which instrument is sounding, when, and in which location, which is critical for answering questions such as "Did the cello on the left play after the drum on the right?". Without such structure, models tend to conflate co-occurring signals or miss temporally offset actions, leading to lower accuracy in complex cross-modal scenarios.

**Spatial-temporal design provides a generalizable advantage across diverse Music AVQA tasks.** Our analysis reveals that models equipped with spatial-temporal design, such as beat-synchronous segment alignment in AMUSE, progressive temporal-spatial filtering in PSTP-NET, and audio-guided token selection in DG-SCT—achieve consistently higher accuracy across audio (e.g., LAST-ATT: 85.71%), visual (e.g., LSTTA: 82.03%), and audio-visual (e.g., AMUSE: 82.43%) question types. These performance gains are particularly pronounced on tasks requiring temporal ordering or cross-modal localization, as shown in Figures 2 and 3. Despite some strong baselines using large-scale pretrained encoders, we observe that models lacking spatial-temporal design struggle with tasks requiring temporal resolution or spatial grounding. Notably, many high-performing models adopt a common architectural pattern: (1) identifying question-relevant time segments, (2) focusing on spatial regions associated with sound cues, and (3) fusing modalities with fine-grained temporal awareness. This recurring design motif underscores spatial-temporal design as not only empirically effective, but also structurally aligned with the demands of reasoning over continuous, densely layered music data.

## 7 A Call on Specialized Musical Designs for Music AVQA

Current Music AVQA models typically treat musical audio as generic acoustic input, operating directly on spectrograms or waveforms without incorporating structured musical attributes such as tempo, downbeats, key, or chord progressions. More fundamentally, human understanding of music relies on hierarchical temporal structure, harmonic organization, and latent causal intent—all of which are shaped by domain-specific knowledge and perceptual priors. Inspired by this observation, we argue that musical audio should not be treated as a raw signal alone, but as a richly structured modality requiring *embed musical priors and inductive structure* into models.

**Incorporating fine-grained musical event cues.** To support precise temporal reasoning over musical events—such as the entrance or exit of specific instruments—models can benefit from auxiliary timestamp supervision derived from musically meaningful proxies. For example, combining waveform peak analysis, Mel-frequency cepstral coefficients (MFCCs), and spectral change detection can help identify dynamic shifts in the audio stream. Beat-tracking algorithms (e.g., from Librosa) can segment audio by rhythm, while pitch-based estimators (e.g., Aubio's YIN) can trace changes in dominant frequency to indicate evolving instrumental activity. These mid-level cues can be used to generate pseudo-labels for training timestamp encoders, enabling models to better localize temporally

anchored events. Embedding such representations into Music AVQA pipelines may improve eventlevel understanding and enhance the interpretability of the model's temporal predictions.

**Embedding mid-level musical structure into multimodal models.** Structured musical features—such as tempo, key, downbeats, and chord progressions—can provide a coherent framework for aligning audio-visual inputs across time. These symbolic or MIR-derived signals offer interpretable, temporally smooth trajectories that reflect the hierarchical organization of music, such as phrases, sections, and transitions. Crucially, they abstract away from low-level waveform fluctuations and offer a musically meaningful scaffold that persists across different genres, tempos, and instrumentation. By integrating them as auxiliary inputs or attention-guiding signals, models may improve their ability to capture long-range dependencies, maintain rhythmic continuity, and resolve ambiguous instrument interactions—especially in polyphonic or ensemble contexts. This structured conditioning can serve as a musical inductive bias, particularly helpful in complex multimodal scenes where overlapping sources challenge simple bottom-up fusion strategies, and where salient events may not be visually or acoustically distinct without temporal alignment cues.

**Modeling latent musical reasoning trajectories.** Many Music AVQA questions require reasoning over implicit causal or temporal relationships—for example, identifying which performer initiated a musical phrase, or determining whether an instrument's entrance shifted the ensemble's dynamic balance. These questions often lack explicit step-level supervision, making it difficult to learn reasoning paths from labels alone. To address this, models can incorporate latent reasoning trajectories: structured internal variables that represent evolving hypotheses about the musical scene. Rather than directly mapping inputs to answers, the model infers intermediate latent states—such as "which instrument is currently leading," "how the rhythmic intensity is changing," or "which performer is preparing to enter"—and updates these states over time as more multimodal evidence arrives. Architecturally, this can be implemented via hierarchical latent-variable models or recurrent variational modules, where latent states encode musical intentions, transitions, or causal flow. These hidden trajectories allow the model to simulate plausible sequences of musical events, enabling it to answer questions that require extrapolating or filling in missing links between observed signals. Crucially, this style of latent reasoning supports robust generalization by embedding inductive structure aligned with how humans infer musical cause and progression—not just surface-level audio-visual co-occurrence.

**Supervising chain-of-thought reasoning in musical QA.** Some musical questions—especially those involving temporal or causal dependencies—require sequential sub-decisions to reach the correct answer. For instance, the question "Which instrument enters after the piano stops?" involves: (1) detecting piano cessation, (2) identifying subsequent onsets, and (3) selecting the earliest new instrument. Rather than treating such questions as black-box classification, models can be explicitly trained to emit intermediate reasoning steps, either through supervised rationales or pseudo-labels derived from MIR-based event detection. This approach—akin to chain-of-thought (CoT) prompting in LLMs—improves transparency, encourages modular subgoal learning, and helps the model maintain alignment across modalities. Moreover, step-wise supervision can highlight failure points in temporal or semantic inference, offering clearer diagnostics for model improvement. In music contexts, CoT chains can incorporate domain-specific steps such as beat alignment, timbre matching, or onset-event attribution. These interpretable intermediate traces not only support higher accuracy on multi-stage queries but also make it easier to identify reasoning shortcuts and dataset biases.

#### 8 Conclusion

This paper underscores the limitations of general-purpose MLLMs for domain-specific tasks such as Music AVQA. Successful Music AVQA necessitates specialized designs tailored to musical content's unique demands: fine-grained audio-visual processing, precise temporal modeling, and integrated domain knowledge. Our central position is that while general models advance, specialized domains require tailored solutions. We call on the research community to: 1) develop more nuanced music understanding benchmarks, and 2) explore hybrid architectures combining MLLM strengths with music-specific components. The future of effective multimodal AI lies not in a universal approach, but in the thoughtful integration of general capabilities with deep, domain-specific expertise, benefiting music understanding and other complex fields.

### References

- [1] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Conference on Computer Vision and Pattern Recognition*, 2023. Available: https://openaccess.thecvf.com/content/CVPR2023/papers/Girdhar\_ImageBind\_ One\_Embedding\_Space\_To\_Bind\_Them\_All\_CVPR\_2023\_paper.pdf.
- [2] Jiaming Han, Kaixiong Gong, Yiyuan Zhang, Jiaqi Wang, Kaipeng Zhang, Dahua Lin, Yu Qiao, Peng Gao, and Xiangyu Yue. Onellm: One framework to align all modalities with language. In *Conference on Computer Vision and Pattern Recognition*, 2024. Available: https://openaccess.thecvf.com/content/CVPR2024/papers/Han\_OneLLM\_One\_ Framework\_to\_Align\_All\_Modalities\_with\_Language\_CVPR\_2024\_paper.pdf.
- [3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. Available: https://arxiv.org/pdf/2502.13923.
- [4] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. arXiv preprint arXiv:2410.21276, 2024. Available: https://arxiv.org/pdf/2410.21276.
- [5] Tian Liang, Jing Huang, Ming Kong, Luyuan Chen, and Qiang Zhu. Querying as prompt: Parameter-efficient learning for multimodal language model. In *Conference on Computer Vision and Pattern Recognition*, 2024. Available: https://openaccess.thecvf. com/content/CVPR2024/papers/Liang\_Querying\_as\_Prompt\_Parameter-Efficient\_ Learning\_for\_Multimodal\_Language\_Model\_CVPR\_2024\_paper.pdf.
- [6] Rongjie Huang, Mingze Li, Dongchao Yang, Jiatong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu, Zhiqing Hong, Jiawei Huang, Jinglin Liu, et al. Audiogpt: Understanding and generating speech, music, sound, and talking head. In AAAI Conference on Artificial Intelligence, 2024. Available: https://arxiv.org/pdf/2304.12995.
- [7] Federico Simonetta, Stavros Ntalampiras, and Federico Avanzini. Multimodal music information processing and retrieval: Survey and future challenges. In *International Workshop on Multilayer Music Representation and Processing*, 2019. Available: https://ieeexplore. ieee.org/stamp/stamp.jsp?tp=&arnumber=8665366.
- [8] Benno Weck, Ilaria Manco, Emmanouil Benetos, Elio Quinton, George Fazekas, and Dmitry Bogdanov. Muchomusic: Evaluating music understanding in multimodal audio-language models. arXiv preprint arXiv:2408.01337, 2024. Available: https://arxiv.org/pdf/2408. 01337.
- [9] Shansong Liu, Atin Sakkeer Hussain, Chenshuo Sun, and Ying Shan. Music understanding llama: Advancing text-to-music generation with question answering and captioning. In *International Conference on Acoustics, Speech and Signal Processing*, 2024. Available: https://arxiv.org/pdf/2308.11276.
- [10] Guangyao Li, Yake Wei, Yapeng Tian, Chenliang Xu, Ji-Rong Wen, and Di Hu. Learning to answer questions in dynamic audio-visual scenarios. In *Conference on Computer Vision and Pattern Recognition*, 2022. Available: https://openaccess.thecvf.com/content/CVPR2022/papers/Li\_Learning\_To\_ Answer\_Questions\_in\_Dynamic\_Audio-Visual\_Scenarios\_CVPR\_2022\_paper.pdf.
- [11] Xiulong Liu, Zhikang Dong, and Peng Zhang. Tackling data bias in musicavqa: Crafting a balanced dataset for unbiased question-answering. In Winter Conference on Applications of Computer Vision, 2024. Available: https: //openaccess.thecvf.com/content/WACV2024/papers/Liu\_Tackling\_Data\_Bias\_ in\_MUSIC-AVQA\_Crafting\_a\_Balanced\_Dataset\_for\_WACV\_2024\_paper.pdf.
- [12] Yuanyuan Jiang and Jianqin Yin. Target-aware spatio-temporal reasoning via answering questions in dynamic audio-visual scenarios. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, 2023. Available: https://aclanthology.org/2023. findings-emnlp.630.pdf.

- [13] Qilang Ye, Zitong Yu, Rui Shao, Xinyu Xie, Philip Torr, and Xiaochun Cao. Cat: Enhancing multimodal large language model to answer questions in dynamic audio-visual scenarios. In *European Conference on Computer Vision*, 2024. Available: https://arxiv.org/pdf/2403. 04640.
- [14] Xingjian Diao, Chunhui Zhang, Tingxuan Wu, Ming Cheng, Zhongyu Ouyang, Weiyi Wu, and Jiang Gui. Learning musical representations for music performance question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, 2024. Available: https://aclanthology.org/2024.findings-emnlp.159.pdf.
- [15] Dinh-Viet-Toan Le, Louis Bigo, Dorien Herremans, and Mikaela Keller. Natural language processing methods for symbolic music generation and information retrieval: a survey. *Computing Surveys*, 2025. Available: https://arxiv.org/pdf/2402.17467.
- [16] Dmitry Bogdanov, Minz Won, Philip Tovstogan, Alastair Porter, and Xavier Serra. The mtg-jamendo dataset for automatic music tagging. In *ICML Machine Learning for Music Discovery Workshop*, 2019. Available: https://github.com/MTG/mtg-jamendo-dataset.
- [17] Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. arXiv preprint arXiv:2311.07919, 2023. Available: https://arxiv.org/pdf/2311.07919.
- [18] Zihao Deng, Yinghao Ma, Yudong Liu, Rongchen Guo, Ge Zhang, Wenhu Chen, Wenhao Huang, and Emmanouil Benetos. Musilingo: Bridging music and text with pre-trained language models for music captioning and query response. *arXiv preprint arXiv:2309.08730*, 2023. Available: https://arxiv.org/pdf/2309.08730.
- [19] Ruibin Yuan, Hanfeng Lin, Yi Wang, Zeyue Tian, Shangda Wu, Tianhao Shen, Ge Zhang, Yuhang Wu, Cong Liu, Ziya Zhou, et al. Chatmusician: Understanding and generating music intrinsically with llm. arXiv preprint arXiv:2402.16153, 2024. Available: https: //arxiv.org/pdf/2402.16153.
- [20] Mengjie Zhao, Zhi Zhong, Zhuoyuan Mao, Shiqi Yang, Wei-Hsiang Liao, Shusuke Takahashi, Hiromi Wakaki, and Yuki Mitsufuji. Openmu: Your swiss army knife for music understanding. arXiv preprint arXiv:2410.15573, 2024. Available: https://arxiv.org/pdf/2410.15573.
- [21] Andrea Agostinelli, Timo I Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, et al. Musiclm: Generating music from text. arXiv preprint arXiv:2301.11325, 2023. Available: https://arxiv.org/ pdf/2301.11325.
- [22] Peiling Lu, Xin Xu, Chenfei Kang, Botao Yu, Chengyi Xing, Xu Tan, and Jiang Bian. Musecoco: Generating symbolic music from text. arXiv preprint arXiv:2306.00110, 2023. Available: https://arxiv.org/pdf/2306.00110.
- [23] Jie Ma, Min Hu, Pinghui Wang, Wangchun Sun, Lingyun Song, Hongbin Pei, Jun Liu, and Youtian Du. Look, listen, and answer: Overcoming biases for audio-visual question answering. In Advances in Neural Information Processing Systems, 2024. Available: https: //openreview.net/pdf?id=twpPD9UMUN.
- [24] Zhangbin Li, Dan Guo, Jinxing Zhou, Jing Zhang, and Meng Wang. Object-aware adaptivepositivity learning for audio-visual question answering. In AAAI Conference on Artificial Intelligence, 2024. Available: https://dl.acm.org/doi/10.1609/aaai.v38i4.28116.
- [25] Guangyao Li, Wenxuan Hou, and Di Hu. Progressive spatio-temporal perception for audiovisual question answering. In *International Conference on Multimedia*, 2023. Available: https://dl.acm.org/doi/pdf/10.1145/3581783.3612293.
- [26] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Information Processing Systems, 2017. Available: https://proceedings.neurips.cc/ paper\_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.

- [27] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *Conference on Computer Vision and Pattern Recognition*, 2022. Available: https://openaccess.thecvf.com/content/CVPR2022/papers/Liu\_Swin\_ Transformer\_V2\_Scaling\_Up\_Capacity\_and\_Resolution\_CVPR\_2022\_paper.pdf.
- [28] Ke Chen, Xingjian Du, Bilei Zhu, Zejun Ma, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Hts-at: A hierarchical token-semantic audio transformer for sound classification and detection. In *International Conference on Acoustics, Speech and Signal Processing*, 2022. Available: https://arxiv.org/pdf/2202.00874.
- [29] Zhifeng Kong, Arushi Goel, Rohan Badlani, Wei Ping, Rafael Valle, and Bryan Catanzaro. Audio flamingo: a novel audio language model with few-shot learning and dialogue abilities. In *International Conference on Machine Learning*, 2024. Available: https://arxiv.org/ pdf/2402.01831.
- [30] Srinivasan Iyer, Xi Victoria Lin, Ramakanth Pasunuru, Todor Mihaylov, Daniel Simig, Ping Yu, Kurt Shuster, Tianlu Wang, Qing Liu, Punit Singh Koura, et al. Opt-iml: Scaling language model instruction meta learning through the lens of generalization. arXiv preprint arXiv:2212.12017, 2022. Available: https://arxiv.org/pdf/2212.12017.
- [31] Benjamin Elizalde, Soham Deshmukh, and Huaming Wang. Natural language supervision for general-purpose audio representations. In *International Conference on Acoustics, Speech and Signal Processing*, 2024. Available: https://arxiv.org/pdf/2309.05767.
- [32] Ying Cheng, Yang Li, Junjie He, and Rui Feng. Mixtures of experts for audio-visual learning. In Advances in Neural Information Processing Systems, 2025. Available: https://openreview. net/pdf?id=SNmuKbU@am.
- [33] Idan Schwartz, Alexander G Schwing, and Tamir Hazan. A simple baseline for audio-visual scene-aware dialog. In *Conference on Computer Vision and Pattern Recognition*, 2019. Available: https://arxiv.org/pdf/1904.05876.
- [34] Yan-Bo Lin and Gedas Bertasius. Siamese vision transformers are scalable audio-visual learners. In European Conference on Computer Vision, 2024. Available: https://www.ecva. net/papers/eccv\_2024/papers\_ECCV/papers/02220.pdf.
- [35] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, 2021. Available: https://proceedings.mlr.press/v139/radford21a/radford21a.pdf.
- [36] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition*, 2016. Available: https://arxiv.org/pdf/1512.03385.
- [37] Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *International Conference on Acoustics, Speech and Signal Processing*, 2017. Available: https://static.googleusercontent.com/media/research. google.com/en//pubs/archive/45857.pdf.
- [38] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023. Available: https://arxiv.org/pdf/2307.09288.
- [39] Zijia Zhao, Longteng Guo, Tongtian Yue, Sihan Chen, Shuai Shao, Xinxin Zhu, Zehuan Yuan, and Jing Liu. Chatbridge: Bridging modalities with large language model as a language catalyst. arXiv preprint arXiv:2305.16103, 2023. Available: https://arxiv.org/pdf/2305.16103.

- [40] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, March 2023. Available: https://lmsys.org/blog/2023-03-30-vicuna/.
- [41] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale. arXiv preprint arXiv:2303.15389, 2023. Available: https://arxiv.org/pdf/2303.15389.
- [42] Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, Wanxiang Che, Xiangzhan Yu, and Furu Wei. Beats: audio pre-training with acoustic tokenizers. In *International Conference on Machine Learning*, 2023. Available: https://arxiv.org/pdf/ 2212.09058.
- [43] Shentong Mo, Weiguo Pian, and Yapeng Tian. Class-incremental grouping network for continual audio-visual learning. In *International Conference on Computer Vision*, 2023. Available: https://openaccess.thecvf.com/content/ICCV2023/papers/Mo\_ Class-Incremental\_Grouping\_Network\_for\_Continual\_Audio-Visual\_Learning\_ ICCV\_2023\_paper.pdf.
- [44] Mingrui Lao, Nan Pu, Yu Liu, Kai He, Erwin M Bakker, and Michael S Lew. Coca: Collaborative causal regularization for audio-visual question answering. In AAAI Conference on Artificial Intelligence, 2023. Available: https://ojs.aaai.org/index.php/AAAI/article/view/ 26527.
- [45] Haytham M. Fayek and Justin Johnson. Temporal reasoning via audio question answering. *Transactions on Audio, Speech, and Language Processing*, 2020. Available: https://arxiv. org/pdf/1911.09655.
- [46] Yuxin Guo, Siyang Sun, Shuailei Ma, Kecheng Zheng, Xiaoyi Bao, Shijie Ma, Wei Zou, and Yun Zheng. Crossmae: Cross-modality masked autoencoders for regionaware audio-visual pre-training. In *Conference on Computer Vision and Pattern Recognition*, 2024. Available: https://openaccess.thecvf.com/content/CVPR2024/ papers/Guo\_CrossMAE\_Cross-Modality\_Masked\_Autoencoders\_for\_Region-Aware\_ Audio-Visual\_Pre-Training\_CVPR\_2024\_paper.pdf.
- [47] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Conference on Computer Vision and Pattern Recognition*, 2022. Available: https://openaccess.thecvf.com/content/ CVPR2022/papers/He\_Masked\_Autoencoders\_Are\_Scalable\_Vision\_Learners\_CVPR\_ 2022\_paper.pdf.
- [48] Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. In *International Conference on Neural Information Processing Systems*, 2022. Available: https://arxiv. org/pdf/2207.06405.
- [49] Changsheng Lv, Shuai Zhang, Yapeng Tian, Mengshi Qi, and Huadong Ma. Disentangled counterfactual learning for physical audiovisual commonsense reasoning. In Advances in Neural Information Processing Systems, 2023. Available: https://proceedings.neurips.cc/paper\_files/paper/2023/file/ 29571f8fda54fe93631c41aad4215abc-Paper-Conference.pdf.
- [50] Pengcheng He, Jianfeng Gao, and Weizhu Chen. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. *arXiv preprint arXiv:2111.09543*, 2021. Available: https://arxiv.org/pdf/2111.09543.
- [51] Yuan Gong, Yu-An Chung, and James Glass. Ast: Audio spectrogram transformer. *arXiv* preprint arXiv:2104.01778, 2021. Available: https://arxiv.org/pdf/2104.01778.
- [52] Haoyi Duan, Yan Xia, Zhou Mingze, Li Tang, Jieming Zhu, and Zhou Zhao. Cross-modal prompts: Adapting large pre-trained models for audio-visual downstream tasks. In Advances in Neural Information Processing Systems, 2023. Available: https://openreview.net/pdf? id=9MwidIH4ea.

- [53] Yaoting Wang, Peiwen Sun, Dongzhan Zhou, Guangyao Li, Honggang Zhang, and Di Hu. Ref-avs: Refer and segment objects in audio-visual scenes. In *European Conference on Computer Vision*, 2024. Available: https://arxiv.org/pdf/2407.10957.
- [54] Liu Zhuang, Lin Wayne, Shi Ya, and Zhao Jun. A robustly optimized BERT pre-training approach with post-training. In *Chinese National Conference on Computational Linguistics*, 2021. Available: https://aclanthology.org/2021.ccl-1.108.pdf.
- [55] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *International Conference on Computer Vision*, 2015. Available: https://openaccess.thecvf.com/content\_iccv\_ 2015/papers/Antol\_VQA\_Visual\_Question\_ICCV\_2015\_paper.pdf.
- [56] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556, 2014. Available: https://arxiv.org/ pdf/1409.1556.
- [57] Thao Minh Le, Vuong Le, Svetha Venkatesh, and Truyen Tran. Hierarchical conditional relation networks for video question answering. In *Conference on Computer Vision and Pattern Recognition*, 2020. Available: https://arxiv.org/pdf/2002.10698.
- [58] Yan-Bo Lin, Yi-Lin Sung, Jie Lei, Mohit Bansal, and Gedas Bertasius. Vision transformers are parameter-efficient audio-visual learners. In *Conference on Computer Vision and Pattern Recognition*, 2023. Available: https://openaccess.thecvf.com/ content/CVPR2023/papers/Lin\_Vision\_Transformers\_Are\_Parameter-Efficient\_ Audio-Visual\_Learners\_CVPR\_2023\_paper.pdf.
- [59] Heeseung Yun, Youngjae Yu, Wonsuk Yang, Kangil Lee, and Gunhee Kim. Pano-avqa: Grounded audio-visual question answering on 360deg videos. In *International Conference on Computer Vision*, 2021. Available: https://arxiv.org/pdf/2110.05122.
- [60] Hongye Liu, Xianhai Xie, Yang Gao, and Zhou Yu. Parameter-efficient transfer learning for audio-visual-language tasks. In *International Conference on Multimedia*, 2023. Available: https://dl.acm.org/doi/pdf/10.1145/3581783.3611939.
- [61] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al. Conformer: Convolution-augmented transformer for speech recognition. *arXiv preprint arXiv:2005.08100*, 2020. Available: https://arxiv.org/pdf/2005.08100.
- [62] Jie Ma, Zhitao Gao, Qi Chai, Jun Liu, Pinghui Wang, Jing Tao, and Zhou Su. Fortisavqa and maven: a benchmark dataset and debiasing framework for robust multimodal reasoning. *arXiv* preprint arXiv:2504.00487, 2025. Available: https://arxiv.org/pdf/2504.00487.
- [63] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Conference on Computer Vision and Pattern Recognition*, 2024. Available: https: //openaccess.thecvf.com/content/CVPR2024/papers/Chen\_InternVL\_Scaling\_ up\_Vision\_Foundation\_Models\_and\_Aligning\_for\_Generic\_CVPR\_2024\_paper.pdf.
- [64] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. Deep modular co-attention networks for visual question answering. In *Conference on Computer Vision and Pattern Recognition*, 2019. Available: https://openaccess.thecvf.com/content\_CVPR\_2019/papers/ Yu\_Deep\_Modular\_Co-Attention\_Networks\_for\_Visual\_Question\_Answering\_CVPR\_ 2019\_paper.pdf.
- [65] Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In Alessandro Moschitti, Bo Pang, and Walter Daelemans, editors, *Conference on Empirical Methods in Natural Language Processing*, 2014. Available: https: //aclanthology.org/D14-1162.pdf.

- [66] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards realtime object detection with region proposal networks. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, 2015. Available: https://proceedings.neurips.cc/paper\_files/paper/2015/file/ 14bfa6bb14875e45bba028a21ed38046-Paper.pdf.
- [67] Sanjoy Chowdhury, Sayan Nag, Subhrajyoti Dasgupta, Jun Chen, Mohamed Elhoseiny, Ruohan Gao, and Dinesh Manocha. Meerkat: Audio-visual large language model for grounding in space and time. In *European Conference on Computer Vision*, 2025. Available: https://www.ecva.net/papers/eccv\_2024/papers\_ECCV/papers/08071.pdf.
- [68] Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. Clap: Learning audio concepts from natural language supervision. In *International Conference on Acoustics, Speech and Signal Processing*, 2023. Available: https://arxiv.org/pdf/2206.04769.
- [69] Yake Wei, Di Hu, Henghui Du, and Ji-Rong Wen. On-the-fly modulation for balanced multimodal learning. *Transactions on Pattern Analysis and Machine Intelligence*, 2024. Available: https://arxiv.org/pdf/2410.11582.
- [70] Xiangpeng Li, Jingkuan Song, Lianli Gao, Xianglong Liu, Wenbing Huang, Xiangnan He, and Chuang Gan. Beyond rnns: positional self-attention with co-attention for video question answering. In AAAI Conference on Artificial Intelligence, 2019. Available: https://doi.org/10.1609/aaai.v33i01.33018658.
- [71] Kunyu Peng, Jia Fu, Kailun Yang, Di Wen, Yufan Chen, Ruiping Liu, Junwei Zheng, Jiaming Zhang, M Saquib Sarfraz, Rainer Stiefelhagen, et al. Referring atomic video action recognition. In *European Conference on Computer Vision*, 2024. Available: https://www.ecva.net/papers/eccv\_2024/papers\_ECCV/papers/02873.pdf.
- [72] Jing Liu, Sihan Chen, Xingjian He, Longteng Guo, Xinxin Zhu, Weining Wang, and Jinhui Tang. Valor: Vision-audio-language omni-perception pretraining model and dataset. *Transactions on Pattern Analysis and Machine Intelligence*, 2025. Available: https: //arxiv.org/pdf/2304.08345.
- [73] Sihan Chen, Handong Li, Qunbo Wang, Zijia Zhao, Mingzhen Sun, Xinxin Zhu, and Jing Liu. Vast: a vision-audio-subtitle-text omni-modality foundation model and dataset. In Advances in Neural Information Processing Systems, 2023. Available: https://proceedings.neurips.cc/paper\_files/paper/2023/file/ e6b2b48b5ed90d07c305932729927781-Paper-Conference.pdf.
- [74] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019. Available: https://aclanthology.org/N19-1423.pdf.
- [75] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. Available: https://openreview.net/pdf?id= YicbFdNTTy.
- [76] Zesen Cheng, Sicong Leng, Hang Zhang, Yifei Xin, Xin Li, Guanzheng Chen, Yongxin Zhu, Wenqi Zhang, Ziyang Luo, Deli Zhao, et al. Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms. arXiv preprint arXiv:2406.07476, 2024. Available: https://arxiv.org/pdf/2406.07476.
- [77] Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen, Meng Zhao, Yifan Zhang, Shaoqi Dong, Xiong Wang, Di Yin, Long Ma, et al. Vita: Towards open-source interactive omni multimodal llm. arXiv preprint arXiv:2408.05211, 2024. Available: https://arxiv.org/pdf/2408. 05211.

- [78] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. arXiv preprint arXiv:2401.04088, 2024. Available: https://arxiv.org/pdf/2401.04088.
- [79] Olga Slizovskaia, Gloria Haro, and Emilia Gómez. Conditioned source separation for musical instrument performances. *Transactions on Audio, Speech, and Language Processing*, 2021. Available: https://arxiv.org/pdf/2004.03873.
- [80] Chuang Gan, Deng Huang, Hang Zhao, Joshua B Tenenbaum, and Antonio Torralba. Music gesture for visual sound separation. In *Conference on Computer Vision and Pattern Recognition*, 2020. Available: https://openaccess.thecvf.com/content\_CVPR\_2020/papers/Gan\_ Music\_Gesture\_for\_Visual\_Sound\_Separation\_CVPR\_2020\_paper.pdf.
- [81] Yitong Jin, Zhiping Qiu, Yi Shi, Shuangpeng Sun, Chongwu Wang, Donghao Pan, Jiachen Zhao, Zhenghao Liang, Yuan Wang, Xiaobing Li, et al. Audio matters too! enhancing markerless motion capture with audio signals for string performance capture. *Transactions on Graphics*, 2024. Available: https://arxiv.org/pdf/2405.04963.
- [82] Sebastian Böck, Filip Korzeniowski, Jan Schlüter, Florian Krebs, and Gerhard Widmer. madmom: A new python audio and music signal processing library. In *International Conference* on *Multimedia*, 2016. Available: https://arxiv.org/pdf/1605.07008.
- [83] Derry Fitzgerald. Harmonic/percussive separation using median filtering. In International Conference on Digital Audio Effects, 2010. Available: https: //www.audiolabs-erlangen.de/media/pages/resources/aps-w23/papers/ 7ef50065d9-1663763115/2010\_FitzGerald\_HarmonicPercussiveSep\_DAFx.pdf.
- [84] Jincheng Huang, Yujie Mo, Ping Hu, Xiaoshuang Shi, Shangbo Yuan, Zeyu Zhang, and Xiaofeng Zhu. Exploring the role of node diversity in directed graph representation learning. In *International Joint Conference on Artificial Intelligence*, 2024. Available: https://www. ijcai.org/proceedings/2024/0229.pdf.
- [85] Peng Zhou, Wei Shi, Jun Tian, Zhenyu Qi, Bingchen Li, Hongwei Hao, and Bo Xu. Attention-based bidirectional long short-term memory networks for relation classification. In Annual Meeting of the Association for Computational Linguistics, 2016. Available: https://aclanthology.org/P16-2034/.
- [86] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. Hierarchical question-image co-attention for visual question answering. In Advances in Neural Information Processing Systems, 2016. Available: https://proceedings.neurips.cc/paper/2016/file/ 9dcb88e0137649590b755372b040afad-Paper.pdf.
- [87] Chenyou Fan, Xiaofan Zhang, Shu Zhang, Wensheng Wang, Chi Zhang, and Heng Huang. Heterogeneous memory enhanced multimodal attention model for video question answering. In *Conference on Computer Vision and Pattern Recognition*, 2019. Available: https://arxiv. org/pdf/1904.04357.
- [88] Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew Zisserman. Vggsound: A large-scale audio-visual dataset. In *International Conference on Acoustics, Speech and Signal Processing*, 2020. Available: https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber= 9053174.
- [89] Yapeng Tian, Jing Shi, Bochen Li, Zhiyao Duan, and Chenliang Xu. Audio-visual event localization in unconstrained videos. In European Conference on Computer Vision, 2018. Available: https://openaccess.thecvf.com/content\_ECCV\_2018/ papers/Yapeng\_Tian\_Audio-Visual\_Event\_Localization\_ECCV\_2018\_paper.pdf#:~: text=In%20particular%2C%20we%20define%20an,what%20category%20the%20event% 20belongs.
- [90] Jinxing Zhou, Jianyuan Wang, Jiayi Zhang, Weixuan Sun, Jing Zhang, Stan Birchfield, Dan Guo, Lingpeng Kong, Meng Wang, and Yiran Zhong. Audio-visual segmentation. In *European Conference on Computer Vision*, 2022. Available: https://arxiv.org/pdf/2207.05042.

- [91] Pinci Yang, Xin Wang, Xuguang Duan, Hong Chen, Runze Hou, Cong Jin, and Wenwu Zhu. Avqa: A dataset for audio-visual question answering on videos. In *International Conference on Multimedia*, 2022. Available: https://dl.acm.org/doi/pdf/10.1145/3503161.3548291.
- [92] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Conference on Computer Vision and Pattern Recognition*, 2016. Available: https://openaccess.thecvf.com/content\_cvpr\_2016/papers/ Xu\_MSR-VTT\_A\_Large\_CVPR\_2016\_paper.pdf.
- [93] Yuncheng Li, Yale Song, Liangliang Cao, Joel Tetreault, Larry Goldberg, Alejandro Jaimes, and Jiebo Luo. Tgif: A new dataset and benchmark on animated gif description. In *Conference on Computer Vision and Pattern Recognition*, 2016. Available: https://openaccess. thecvf.com/content\_cvpr\_2016/papers/Li\_TGIF\_A\_New\_CVPR\_2016\_paper.pdf.
- [94] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *International Conference on Computer Vision*, 2017. Available: https://openaccess.thecvf.com/content\_ICCV\_2017/papers/ Krishna\_Dense-Captioning\_Events\_in\_ICCV\_2017\_paper.pdf.
- [95] Ramon Sanabria, Ozan Caglayan, Shruti Palaskar, Desmond Elliott, Loïc Barrault, Lucia Specia, and Florian Metze. How2: a large-scale dataset for multimodal language understanding. arXiv preprint arXiv:1811.00347, 2018. Available: https://arxiv.org/pdf/1811.00347.
- [96] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. Vatex: A large-scale, high-quality multilingual dataset for video-and-language research. In *International Conference on Computer Vision*, 2019. Available: https://openaccess. thecvf.com/content\_ICCV\_2019/papers/Wang\_VaTeX\_A\_Large-Scale\_High-Quality\_ Multilingual\_Dataset\_for\_Video-and-Language\_Research\_ICCV\_2019\_paper.pdf.
- [97] Anthony Colas, Seokhwan Kim, Franck Dernoncourt, Siddhesh Gupte, Zhe Wang, and Doo Soon Kim. TutorialVQA: Question answering dataset for tutorial videos. In *Language Resources and Evaluation Conference*, 2020. Available: https://aclanthology.org/2020. lrec-1.670.pdf.
- [98] Shaojie Wang, Wentian Zhao, Ziyi Kou, Jing Shi, and Chenliang Xu. How to make a blt sandwich? learning vqa towards understanding web instructional videos. In Winter Conference on Applications of Computer Vision, 2021. Available: https://openaccess.thecvf.com/content/WACV2021/papers/Wang\_How\_to\_Make\_a\_ BLT\_Sandwich\_Learning\_VQA\_Towards\_Understanding\_WACV\_2021\_paper.pdf.
- [99] Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. Just ask: Learning to answer questions from millions of narrated videos. In *International Conference* on Computer Vision, 2021. Available: https://openaccess.thecvf.com/content/ ICCV2021/papers/Yang\_Just\_Ask\_Learning\_To\_Answer\_Questions\_From\_Millions\_ of\_Narrated\_ICCV\_2021\_paper.pdf.
- [100] Huda Alamri, Vincent Cartillier, Abhishek Das, Jue Wang, Anoop Cherian, Irfan Essa, Dhruv Batra, Tim K Marks, Chiori Hori, Peter Anderson, et al. Audio visual scene-aware dialog. In *Conference on Computer Vision and Pattern Recognition*, 2019. Available: https://openaccess.thecvf.com/content\_CVPR\_2019/papers/ Alamri\_Audio\_Visual\_Scene-Aware\_Dialog\_CVPR\_2019\_paper.pdf.
- [101] Huda Alamri, Vincent Cartillier, Raphael Gontijo Lopes, Abhishek Das, Jue Wang, Irfan Essa, Dhruv Batra, Devi Parikh, Anoop Cherian, Tim K Marks, et al. Audio visual sceneaware dialog (avsd) challenge at dstc7. arXiv preprint arXiv:1806.00525, 2018. Available: https://arxiv.org/pdf/1806.00525.
- [102] Luis Fernando D'Haro, Koichiro Yoshino, Chiori Hori, Tim K Marks, Lazaros Polymenakos, Jonathan K Kummerfeld, Michel Galley, and Xiang Gao. Overview of the seventh dialog system technology challenge: Dstc7. *Computer Speech & Language*, 2020. Available: https://jkk.name/pub/csl20dstc.pdf.

- [103] Ankit Shah, Shijie Geng, Peng Gao, Anoop Cherian, Takaaki Hori, Tim K Marks, Jonathan Le Roux, and Chiori Hori. Audio-visual scene-aware dialog and reasoning using audio-visual transformers with joint student-teacher learning. In *International Conference on Acoustics*, *Speech and Signal Processing*, 2022. Available: https://ieeexplore.ieee.org/stamp/ stamp.jsp?tp=&arnumber=9746481.
- [104] Amir Zadeh, Rowan Zellers, Eli Pincus, and Louis-Philippe Morency. Mosi: multimodal corpus of sentiment intensity and subjectivity analysis in online opinion videos. arXiv preprint arXiv:1606.06259, 2016. Available: https://arxiv.org/pdf/1606.06259.
- [105] Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J. Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences. In Annual Meeting of the Association for Computational Linguistics, 2019. Available: https://aclanthology.org/P19-1656.pdf.
- [106] Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N Chang, Sungbok Lee, and Shrikanth S Narayanan. Iemocap: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 2008. Available: https://link.springer.com/content/pdf/10.1007/s10579-008-9076-6.pdf.
- [107] Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. MELD: A multimodal multi-party dataset for emotion recognition in conversations. In Annual Meeting of the Association for Computational Linguistics, 2019. Available: https://aclanthology.org/P19-1050.pdf.
- [108] Wenmeng Yu, Hua Xu, Fanyang Meng, Yilin Zhu, Yixiao Ma, Jiele Wu, Jiyun Zou, and Kaicheng Yang. CH-SIMS: A Chinese multimodal sentiment analysis dataset with finegrained annotation of modality. In Annual Meeting of the Association for Computational Linguistics, 2020. Available: https://aclanthology.org/2020.acl-main.343.pdf.
- [109] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, 2021. Available: https://proceedings.mlr.press/v139/jia21b/jia21b.pdf.
- [110] Xiaohua Zhai, Xiao Wang, Basil Mustafa, Andreas Steiner, Daniel Keysers, Alexander Kolesnikov, and Lucas Beyer. Lit: Zero-shot transfer with locked-image text tuning. In *Conference on Computer Vision and Pattern Recognition*, 2022. Available: https://openaccess.thecvf.com/content/CVPR2022/papers/Zhai\_LiT\_ Zero-Shot\_Transfer\_With\_Locked-Image\_Text\_Tuning\_CVPR\_2022\_paper.pdf.
- [111] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. In Advances in Neural Information Processing Systems, 2023. Available: https://proceedings.neurips.cc/paper\_ files/paper/2023/file/90ce332aff156b910b002ce4e6880dec-Paper-Datasets\_and\_ Benchmarks.pdf.
- [112] Binzhu Xie, Sicheng Zhang, Zitang Zhou, Bo Li, Yuanhan Zhang, Jack Hessel, Jingkang Yang, and Ziwei Liu. Funqa: Towards surprising video comprehension. In *European Conference on Computer Vision*, 2024. Available: https://arxiv.org/pdf/2306.14899.

# **Contents of Appendix**

A	Quantitative Comparison on Music AVQA Datasets	20					
B	Representative Examples of Music Performance Scene Types         21						
С	C Representative Examples of Music AVQA Question Types21						
D	How Music AVQA Differs from Traditional Multimodal Tasks	23					
E	Details of Music AVQA Datasets	24					
	E.1 Music AVQA Datasets	24					
	E.2 Key Difference from Other AVQA Datasets	25					
F	Details of Music AVQA Methods with Spatial-Temporal Designs	26					
G	Details of Existing Music AVQA Methods	28					

# A Quantitative Comparison on Music AVQA Datasets

We present comprehensive quantitative comparisons of recent state-of-the-art methods on multiple Music AVQA datasets [10, 11, 23], shown in Table 2, 3, and 4. We evaluate the models across a diverse set of question categories, spanning Audio-related, Visual-related, and Audio&Visual-related reasoning tasks. For each dataset, we report accuracy metrics for subcategories such as Counting, Comparative, Location, Existential, and Temporal reasoning, along with average accuracy within each modality and the overall performance.

Table 2: Comparison with state-of-the-art methods on the Music AVQA [10] test set. We report the accuracy for Audio (Counting, Comparative), Visual (Counting, Location), and Audio-Visual (Existential, Counting, Location, Comparative, Temporal) question types, along with the average accuracy for Audio, Visual, Audio-Visual, and overall. Methods highlighted with a gray background incorporate spatial-temporal designs.

Methods	Aud	lio-related	I QA	Visu	al-relate	d QA	Audio&Visual-related QA						Ανσ
memous	Count	Comp	Avg	Count	Local	Avg	Exist	Count	Local	Comp	Temp	Avg	n's
Amuse [14]	84.61	82.45	83.58	87.41	84.39	85.84	86.95	85.49	73.01	82.98	83.06	82.43	83.52
AUDIO FLAMINGO [29]	-	-	-	-	-	-	-	-	-	-	-	-	-
AVMOE [32]	-	-	77.60	-	-	82.70	-	-	-	-	-	71.90	75.70
AVSD [33]	72.41	61.90	68.52	67.39	74.19	70.83	81.61	58.79	63.89	61.52	61.41	65.49	67.44
AVSIAM [34]	-	-	-	-	-	-	-	-	-	-	-	-	-
AVST [10]	77.78	67.17	73.87	73.52	75.27	74.40	82.49	69.88	64.24	64.67	65.82	69.53	71.59
CAT [13]	-	-	84.90	-	-	86.10	-	-	-	-	-	83.20	84.30
CHATBRIDGE [39]	-	-	28.90	-	-	33.10	-	-	-	-	-	43.00	78.90
CIGN [43]	-	-	-	-	-	-	-	-	-	-	-	-	-
COCA [44]	79.94	67.68	75.42	75.10	75.43	75.23	83.50	66.63	69.72	64.12	65.57	69.96	72.33
CONVLSTM [45]	68.88	63.06	66.73	64.89	58.55	61.68	82.81	55.99	61.30	53.45	54.73	61.75	62.61
CROSSMAE [46]	-	-	-	-	-	-	-	-	-	-	-	-	-
DCL [49]	-	-	-	-	-	-	-	-	-	-	-	-	-
DG-SCT [52]	83.27	64.56	76.34	81.57	82.57	82.08	81.61	72.84	65.91	64.22	67.48	70.56	74.62
EEMC [53]	-	-	-	-	-	-	-	-	-	-	-	-	-
FCNLSTM [45]	69.96	61.06	66.67	63.89	58.14	60.98	83.42	56.31	60.28	50.85	56.92	61.46	62.25
GPT-40 [4]	65.42	36.07	50.75	72.36	62.30	67.33	56.12	54.84	59.23	37.84	42.35	50.08	54.06
GRU [55]	71.82	58.90	67.04	66.06	71.82	68.97	81.41	60.30	62.32	56.23	61.89	64.26	66.00
HCRN [57]	70.21	45.62	61.14	62.41	51.51	56.90	52.94	42.07	54.70	50.59	33.33	48.41	52.54
LAST-ATT [11]	85.71	63.10	-	83.86	83.09	-	76.47	76.20	68.91	65.60	66.75	-	75.45
LAVISH [58]	75.59	84.13	76.86	77.45	72.91	76.29	71.91	77.52	75.81	76.75	77.62	76.31	76.10
LAVIT [59]	74.36	64.56	70.73	69.39	75.65	72.56	81.21	59.33	64.91	64.22	63.23	66.64	68.93
LSTTA [60]	81.75	82.04	81.90	81.82	82.23	82.03	83.46	79.11	78.23	78.02	79.32	79.63	81.19
MAVEN [62]	79.44	54.10	72.79	80.49	93.50	86.99	87.00	66.67	73.85	54.95	68.24	69.94	74.60
MCAN [64]	/5.05	54.58	67.47	68.06	72.15	70.13	81.91	54.15	53.45	52.11	47.21	57.80	62.77
MCCD [23]	83.87	/1.04	/9.14	/9./8	/6./3	/8.24	80.87	51.63	/1.46	64.67	64.60	67.13	72.20
MEERKAT [6/]	-	-	-	-	-	-	-	85.70	-	/5.98	-	-	-
ONELLM [2]	-	-	-	-	-	-	-	-	-	-	-	-	47.60
OPM [69]	-	-	-	-	-	-	-	-	-	-	-	-	/0.80
PSAC [70]	/1.33	56.07	05.08	05.89	72.07	69.02	76.59	54.80	03.11	55.96	61.17	62.75	64.92
PSTP-NET [25]	13.91	65.59	70.90	//.15	//.30	//.20	/0.18	13.23	/1.80	/1.19	69.00	12.57	13.52
QAP [5]	-	-	-	-	-	-	-	-	-	-	-	-	-
QWEN2.3-VL [3]	48.00	55.00	51.80	33.28	35.00	34.47	44.00	32.17	05.57	37.64	41.10	47.75	50.14
KEFAIOMINET [/1]	-	-	-	-	-	-	-	-	-	-	-	-	- 78.00
VALUK [/2]	-	-	74.06	- 71.56	-	74.20	- 01 01	- 64.51	-	-	-	13.30 60.54	71.50
VIDEOLIAMA 2 [76]	70.10	52.46	60.64	×1.30	20.36	22.11	77.00	62.44	77.60	50.46	64 71	68.08	72.56
VIDEOLLAWA-2 [70]	79.44 50.81	45.00	54.76	50.41	34.06	02.11	54.00	40.46	46.02	27.02	41.19	12 74	12.50
VIIA [//]	39.81	45.90	54.70	50.41	54.90	42.00	54.00	49.40	40.92	21.95	41.10	45.74	45.44

Table 3: Comparison with state-of-the-art methods on the Music AVQA v2.0 [11] test set. We report the accuracy for Audio (Counting, Comparative), Visual (Counting, Location), and Audio-Visual (Existential, Counting, Location, Comparative, Temporal) question types, along with the average accuracy for Audio, Visual, Audio-Visual, and overall. Methods highlighted with a gray background incorporate spatial-temporal designs.

	r -		0										
Methods	Aud	Audio-related QA		Visual-related QA			Audio&Visual-related QA						Avg
	Count	Comp	Avg	Count	Local	Avg	Exist	Count	Local	Comp	Temp	Avg	
Amuse [14]	84.76	83.88	84.34	88.15	85.16	86.74	88.30	87.47	78.77	84.41	85.38	85.51	85.16
AVST [10]	81.74	62.11	72.46	79.08	77.64	78.40	72.12	69.03	65.05	63.98	60.57	66.26	71.08
DG-SCT [52]	83.66	62.47	73.64	82.05	82.97	82.48	83.43	72.70	64.65	64.78	67.34	70.38	74.08
LAST-ATT [11]	86.03	62.52	-	84.12	84.01	-	76.21	75.23	68.91	65.60	60.60	-	75.44
LAVISH [58]	84.36	58.57	72.17	83.25	81.46	82.40	73.26	73.45	65.64	64.26	60.82	67.75	72.34

Table 4: Comparison with state-of-the-art methods on the Music-AVQA-R [23] test set. We report the accuracy for Audio (Counting, Comparative), Visual (Counting, Location), and Audio-Visual (Existential, Counting, Location, Comparative, Temporal) question types, along with the average accuracy for Audio, Visual, Audio-Visual, and overall. Methods highlighted with a gray background incorporate spatial-temporal designs.

Methods	Aud	Audio-related QA		Visual-related QA			Audio&Visual-related QA						Δνσ
	Count	Comp	Avg	Count	Local	Avg	Exist	Count	Local	Comp	Temp	Avg	
ATT-BLSTM [85]	60.00	49.55	54.77	32.15	47.97	48.89	54.33	39.46	32.52	51.00	24.45	40.35	40.35
AVSD [33]	50.92	54.20	52.56	35.21	68.11	52.20	64.13	36.68	27.14	58.99	40.83	45.55	45.55
CONVLSTM [45]	55.68	60.22	57.95	35.64	51.66	52.23	72.45	53.18	32.35	57.91	43.33	51.84	51.84
FCNLSTM [45]	51.36	57.96	54.66	33.53	52.96	50.09	71.64	51.98	34.96	57.40	33.90	49.98	49.98
GRU [55]	57.78	58.95	58.36	38.08	57.67	54.17	70.53	43.33	39.70	57.29	35.85	49.34	49.34
HCATTN [86]	51.65	53.12	52.38	32.86	60.09	50.02	63.85	39.77	36.01	54.47	36.54	46.13	46.13
HCRN [57]	54.42	39.81	47.11	32.71	45.34	43.88	53.63	39.67	37.08	35.10	42.30	41.56	41.56
HME [87]	58.28	56.63	57.45	33.71	65.93	54.40	66.12	39.91	40.18	56.89	37.76	48.17	48.17
LAVISH [58]	52.86	62.72	57.79	38.33	67.47	55.83	78.65	41.48	32.38	62.18	44.05	51.75	51.75
LAVIT [59]	47.01	47.86	47.43	31.39	66.35	48.01	37.21	53.02	36.87	43.05	42.17	42.46	42.46
MCAN [64]	67.59	54.49	61.04	45.64	64.37	58.62	59.29	53.86	45.02	51.49	46.35	51.20	51.20
MCCD [23]	75.78	63.43	69.60	61.76	73.43	68.80	76.18	50.55	50.92	62.15	66.95	61.35	61.35
PSAC [70]	54.85	52.77	53.81	37.99	66.83	53.25	53.05	47.14	38.14	48.53	36.46	44.66	44.66

## **B** Representative Examples of Music Performance Scene Types

Figure 4 presents examples for each of the scene types defined in Section 2. These examples further underscore the performance diversity that Music AVQA methods must accommodate, ranging from sparse solo stages to densely populated, culturally nuanced ensembles.



Figure 4: Representative examples for the four common music performance scene types. (a) *Solo performance*: a single musician highlights individual virtuosity on one instrument. (b) *Ensemble of the same instrument*: multiple players of identical (or closely related) instruments create timbral thickness and homogeneous harmony. (c) *Ensemble of different instruments*: a heterogeneous group blends distinct tonal colours and enables richer contrapuntal interaction. (d) *Culture-specific ensemble*: a traditional instrumental configuration (e.g. guzheng quartet, gamelan group) that captures the performance idioms of a particular musical culture.

## C Representative Examples of Music AVQA Question Types

Table 5 and Figure 5 provide representative examples of the Music AVQA question types. These illustrate how each type manifests itself across audio, visual and audio-visual modalities, highlighting the multimodal and fine-grained nature of the task.

Modality	Task Type	Question	Answer
Audio	Counting	Are there acoustic guitar and accordion sound?	Yes
	Comparative	Is the clarinet playing longer than the drum?	No
Visual	Counting	Are there violin and ukulele instruments in the video?	Yes
	Localization	What kind of musical instrument is it?	Cello
Audio Visual	Existential	Is there a voiceover?	Yes
	Counting	How many instruments are sounding in the video?	Three
	Localization	Is the first sound coming from the middle instrument?	Yes
	Comparative	Is the tuba on the right more rhythmic than the piano?	Yes
	Temporal	Which instrument makes sounds before the violin?	Cello

Table 5: Examples of questions in Music AVQA categorized by modality involved and task type.



Figure 5: Examples of Music AVQA question types spanning audio, visual, and audio-visual modalities, including counting, comparison, localization, existential, and temporal QA.

## D How Music AVQA Differs from Traditional Multimodal Tasks

Music AVQA is unique among audio-visual-text (AVT) tasks because it forces the model to disentangle *polyphonic*, *continuous* sound streams, bind them to precise visual sources, and reason with explicit musical knowledge.

**Generic Audio–Visual Question Answering.** Classic AVQA benchmarks [37, 88, 89, 90] present short clips with a single audible event (e.g., a dog bark) and ask coarse "what/where/when" questions that can be answered once the sounding object is localized and its time span identified. The AVQA dataset itself illustrates this design: most videos contain one foreground sound and minimal polyphony, so millisecond-level alignment is unnecessary [10, 91]. Music AVQA, in contrast, poses queries such as "Which violin enters after the flute?", which requires the model to track multiple *overlapping* instruments and reason over precise temporal order, a level of granularity generic AVQA never targets.

**Video Captioning.** Datasets like MSR-VTT [92], GIF [93], ActivityNet Captions [94], How2 [95], and Vatex [96] evaluate whether a system can produce one or two fluent sentences that summarize the *gist* of a clip. Moreover, temporal mis-alignment of a few seconds or the omission of background sounds rarely affects the score. Music AVQA eliminates summarisation entirely and replaces it with beat-accurate, question-driven reasoning: the model must pinpoint note onsets, match them to performers, and compare rhythmic or dynamic patterns tasks far beyond the scope of generic captioning.

**Instructional Video QA.** TutorialVQA [97], YouCookQA [98] and HowToVQA69M [99] frame question answering around narrated procedural videos whose audio channel is dominated by speech that explicitly describes each step. The narration acts as a guide-track and overlapping non-speech sounds are rare and not queried. Music AVQA removes narration altogether and treats the dense musical audio as the primary reasoning target, forcing models to infer structure (beats, phrases, instrument entrances) directly from raw sound rather than from textual guidance.

Audio–Visual Scene-Aware Dialogue. The AVSD dataset [100] and its follow-up [101, 33, 102, 103] train systems to hold multi-turn conversations about short ( $\approx$  10 s) household videos, grounding answers in coarse scene context while maintaining dialog coherence. Acoustic events are typically brief (speech, clatter) and alignment at  $\pm 1$  s suffices. Music AVQA disregards conversational flow and instead demands hierarchical timing: beats, bars, sections. Every answer depends on tight audio–video synchrony, not on dialog history management.

**Multimodal Sentiment Analysis.** Benchmarks such as CMU-MOSI [104] and CMU-MOSEI [105] fuse text transcripts, facial expression and prosody to predict sentiment or emotion over 5 to 30 second clips. Overlap between speakers or sound sources is noise to suppress, and no explicit sourcebinding is required. In Music AVQA overlap is the *signal*: models must isolate each instrument's contribution, bind it to its on-screen performer and reason about their relationships (loudness, order, count). Other emotion datasets, like IEMOCAP [106], MELD [107], and CH-SIMS [108], use 5–30 s clips, treat overlapping voices as background noise, and label sentiment at the utterance—rather than source—level.

**Cross-Modal Retrieval.** Contrastive systems such as CLIP [35] (vision–text), ALIGN [109], LiT [110], or ImageBind [1] (audio–vision–text) learn global embeddings and judge success by top-k similarity—minor temporal or spatial errors barely change the score. Music AVQA instead penalises any mis-alignment: swapping left/right instruments or missing a single beat flips the answer. The task therefore demands persistent token-level grounding rather than coarse embedding proximity.

Where existing AVT tasks rely on sparse events, speech cues or global embeddings, Music AVQA alone couples *dense polyphonic audio*, *frame-accurate audio–visual alignment*, and *music-theoretic knowledge*. It thus sets a much higher bar for multimodal reasoning than traditional benchmarks.

# E Details of Music AVQA Datasets

We summarize the three publicly released Music-AVQA benchmarks and their successive extensions (§ E.1), then contrasts Music-AVQA with several general-purpose AVQA datasets to underline the domain-specific challenges posed by musical performance videos (§ E.2).

## E.1 Music AVQA Datasets

Table 6 provides a summary of three representative datasets specifically designed for Music Performance Audio-Visual Question Answering (Music AVQA) tasks.

Table 6: Evolution and characteristics of Music AVQA datasets: a comparative overview of MUSIC-AVQA [10], MUSIC-AVQA v2.0 [11], and MUSIC-AVQA-R [23] benchmarks.

Dataset	Brief Description
MUSIC-AVQA [10]	The MUSIC-AVQA dataset represents a significant contribution to audio-visual question answering research, comprising 9,288 videos with over 150 hours of musical performances covering 22 instruments, generating 45,867 question-answer pairs. The dataset is randomly split into training, validation, and testing sets with 32,087, 4,595, and 9,185 QA pairs respectively, spanning 33 question templates across 9 question types including existential, location, counting, comparative, and temporal questions.
MUSIC-AVQA v2.0 [11]	The MUSIC-AVQA v2.0 dataset builds upon the original MUSIC-AVQA by ad- dressing data bias issues, comprising 10,518 videos (9,288 from the original plus 1,230 new videos) with musical performances covering 22 instruments, generat- ing approximately 54,000 question-answer pairs. The balanced dataset splits into training, validation, and testing sets with 36,700, 5,250, and 10,819 QA pairs respectively, spanning 33 question templates across 9 question types. The authors specifically balance 15 biased templates by ensuring no dominant answers exceed 60% for binary questions or 50% for multi-class questions, particularly enhancing representation of underrepresented answers in existential, counting, temporal, loca- tion, and comparative question categories.
MUSIC-AVQA-R [23]	The MUSIC-AVQA-R dataset proposed in this paper is an extension of MUSIC-AVQA specifically designed to evaluate the robustness of audio-visual question answering models. It expands the original test set through a human-machine collaboration mechanism that rephrases each question 25 times, increasing the number of questions from 9,129 to 211,572, and introduces distribution shifts to categorize questions into head (common) and tail (rare) samples. Compared to the original dataset, MUSIC-AVQA-R features a vocabulary size of 465 (five times larger than MUSIC-AVQA), provides more diverse question formulations while preserving inherent biases in the training and validation sets, and offers three evaluation metrics—head accuracy, tail accuracy, and overall accuracy—enabling researchers to assess model performance in both in-distribution and out-of-distribution scenarios, making it the first dataset specifically designed for robustness evaluation in audio-visual question answering tasks.

## E.2 Key Difference from Other AVQA Datasets

Table 7 contrasts the Music-AVQA dataset [10] with several widely-used AVQA benchmarks [91, 111, 112, 72, 88]. For each dataset, we highlight the most salient divergence from the music-specific setting, focusing on aspects such as task format, content domain, temporal scope, and the presence or absence of fine-grained musical reasoning.

Table 7: Other representative benchmarks (AVQA [91], EgoSchema [111], FunQA [112], VALOR-1M [72], and VGG-Sound [88]) and the key difference each bears with respect to Music-AVQA [10].

Dataset	Key Difference
AVQA [91]	Builds multiple-choice QA on everyday VGG-Sound clips; questions target generic activities and causal relations in real-life videos, so it lacks the fine-grained instrument/sound localization and music-theory knowledge required in MUSIC-AVQA.
EgoSchema [111]	Uses first-person (Ego4D) footage that is three-minutes long, stressing long-range temporal reasoning in egocentric daily tasks; audio cues are incidental and the task is 5-way multiple choice, very different from the short, professionally filmed music performances and open-ended answers in MUSIC-AVQA.
FunQA [112]	Focuses on "surprising" humour/creative/magic clips (4.3 k videos, 312 k QAs) that test commonsense violations; audio is often background and questions centre on counter-intuitive visual events, not on synchronised musical notes or instrument semantics.
VALOR-1M [72]	A pre-training corpus (1 M videos) with tri-modal captions meant for retrieval/cap- tioning; QA supervision is extremely sparse and relies on auto-generated subtitles, so it serves as a foundation model resource rather than a targeted AVQA evaluation set like MUSIC-AVQA.
VGG-Sound [88]	It is an audio-visual correspondent dataset consisting of short clips of audio sounds from YouTuBe. And it provides raw audio–visual correspondence but no ques- tion–answer supervision or fine-grained reasoning labels.

# F Details of Music AVQA Methods with Spatial-Temporal Designs

Table 8 illustrates methods incorporating explicit **spatial-temporal** design components in detailed.

Table 8: Description of representative methods for spatial-temporal design for Music AVQA.

Paper/Work	Brief Description
Amuse [14]	Focuses on music performance scenarios by aligning time segments in both audio and video streams via a cross-attention paradigm. Exploits synchronized features (such as beat-level or note-level alignment) to capture subtle temporal dependencies among instruments in dense music passages. By integrating rhythmic cues and cross-modal interactions, it is particularly suited for questions that involve multiple instruments playing simultaneously or changing their patterns over time.
AVST [10]	Proposes a spatio-temporal grounded audio-visual approach that explicitly localizes sounding objects in each frame while applying a question-guided temporal attention mechanism. The model grounds audio-visual events and emphasizes which frames (visual) and which segments (audio) are most relevant for question answering. By combining localized visual features and temporal cues, it captures object interactions over time and can better handle questions involving spatial and temporal reasoning.
CIGN [43]	Learns audio-visual class tokens and an Audio-Visual Continual Grouping module that, at every time-step, pulls together frame-level spectrogram features and region features into category-aware clusters. A token-distillation schedule preserves past knowledge while the regrouping logic tracks objects and sounds through the video's timeline, giving the model temporally consistent, cross-modal semantics for spatial-temporal reasoning.
DCL [49]	Introduces a Disentangled Counterfactual Learning framework to handle physical audio-visual commonsense reasoning tasks. Decomposes video signals into static (time-invariant) and dynamic (time-varying) factors using a VAE-based encoder, enabling clearer separation of constant background features from changing events. Additionally employs a counterfactual intervention module on the dynamic factors to perform causal reasoning, helping the model answer "what if" questions related to temporal and event relationships.
DG-SCT [52]	Introduces a Dual-Guided Spatial-Channel-Temporal (DG-SCT) attention layer that is injected in every frozen audio and visual transformer block. Audio prompts steer visual tokens (and vice-versa) via bidirectional attention that highlights salient spatial regions, discriminative feature channels, and pivotal temporal segments, producing fine-grained spatio-temporal alignments that boost related tasks.
EEMC [53]	Divides audio/video into 1-s slices and fuses them with text through a Temporal Bi-modal Transformer backed by a cached-memory mechanism that magnifies sudden changes across time. The resulting multimodal cue stream then serves as a cross-attention prompt for the segmentation decoder, enabling precise localisation of objects and events as their spatial footprints and temporal order evolve.
LAST-ATT [11]	Tackles audio-visual question answering with a repeated cross-attention strategy. Uses Swin-Transformer-v2 for visual frame features and a specialized Audio Spectrogram Transformer for audio, then merges them based on the question. By repeatedly "attending" to the most relevant frames and spectrogram patches, it effectively localizes musical actions (e.g., a pianist's keystrokes) over time. This design is well suited for intricate temporal queries and locating key audio events in dense musical content.

Continued on next page

Paper / Work	Brief Description
LAVISH [58]	Adds a lightweight Latent Audio-Visual HYbrid adapter to every layer of a
	frozen ViT. A compact pool of latent tokens acts as a cross-attention
	bottleneck, letting audio frames gate visual tokens (and vice-versa) as the
	video unfolds, so spatial patches and framewise dynamics are fused early
	while keeping the backbone frozen.
LAVIT [59]	Targets 360° videos with a transformer that augments each patch by a
	quaternion-based spherical coordinate and aligns it with audio via joint
	contrastive objectives. The spherical embedding plus an auxiliary
	audio-skewness prediction head lets the model reason about where (on the
	sphere) and when a sound arises, delivering fine-grained spatial-temporal
	grounding beyond normal FOV clips.
LSTTA [60]	A parameter-efficient transfer learning approach for audio-visual-language
	tasks that adds dedicated adapter modules while freezing large pretrained
	backbones. Splits temporal modeling into two scales: a short-term
	semantic interaction module (for capturing local correlations such as brief
	instrumental flourishes) and a long-term semantic filtering module (for
	broader progressions over many frames). This structure helps the model
	identify when, how, and for how long different instruments contribute,
	achieving a refined spatio-temporal representation.
MAVEN [62]	Employs a Multimodal Audio-Visual Epistemic Network that cycles
	between audio, video and text logits, using debiasing constraints to keep
	modality-specific and fused predictions consistent over time. The cycle
	guidance implicitly anchors each question to the correct temporal
	segments while suppressing spurious correlations.
MCCD [23]	Introduces a Multifaceted Cycle-Collaborative Debiasing objective: KL
	penalties enlarge the gap between uni-modal and tri-modal logits at every
	timestep, then force the three unimodal paths to agree with each other.
	This temporal-cycle training steers attention toward frames (and sounds)
	that all modalities truly share, yielding stabler spatial-temporal grounding
	under distribution shift.
MEERKAT [6/]	Employs a two-stage mechanism for fine-grained audio-visual grounding
	in space and time. First uses an Audio-visual Optimal Transport (AvOp1)
	module for fine-grained local alignment between audio features and
	Specific image patches. Next, the Audio-Visual Attention Consistency
	Enforcement (AVACE) module remes cross-modal attention maps to
	precisely locate audio sources within bounding boxes, enforcing spatial
	that correspond to the audio signal
DSTD NET [25]	Proposes a Progressive Spatio Temporal Perception framework for
1511-NEI [25]	audio-visual $\Omega A$ Divides the selection of relevant information into three
	modules: (1) the Temporal Segment Selection Module (TSSM) for picking
	key time segments pertinent to the question: (2) the Spatial Region
	Selection Module (SRSM) to identify essential visual patches within those
	segments: and (3) the Audio-guided Visual Attention Module (AVAM) to
	align selected visual patches with the audio signals. This stepwise process
	helps isolate question-relevant data and reduce interference.
REFATOMNET [71]	For referring atomic actions, it runs three streams—visual, text and
[,-]	location-semantic tokens— and merges them through novel cross-stream
	agent-attention blocks. The location-semantic stream provides per-person
	bounding-box hints over time, letting the network lock onto the described
	individual before classifying frame-level atomic actions, thus tightly
	coupling spatial localisation with temporal action cues.

Continued on next page

Paper / Work	Brief Description
VIDEOLLAMA-2 [76]	Builds a video-LLM around a Spatial-Temporal Convolution (STC)
	connector that first performs per-frame spatial mixing and then
	downsamples temporally, giving the language model a compact yet
	order-aware token sequence. A jointly-trained audio branch injects
	synchronized spectrogram tokens, enabling the model to answer
	audio-visual questions that hinge on both where events happen on screen
	and when they unfold.

## G Details of Existing Music AVQA Methods

- AVMOE [32]: The paper proposes a parameter-efficient transfer learning framework for audiovisual tasks by dynamically integrating intra-modal and inter-modal information through a mixture of experts. The approach introduces unimodal adapters to capture within-modality details and cross-modal adapters to model interactions between audio-visual streams, while a lightweight modality-agnostic router dynamically allocates expert weights based on input characteristics. By combining these components, AVMOE adaptively balances modality-specific and cross-modal features, addressing challenges like missing modalities or noisy inputs, thereby enhancing robustness and performance across diverse audio-visual tasks such as AV localization, segmentation, and question answering without requiring full model fine-tuning.
- AVSD [33]: The paper proposes an end-to-end baseline for audio-visual scene-aware dialog to
  enhance virtual assistants by integrating multimodal signals. The method employs an attention
  mechanism to differentiate useful signals from distractions, while maintaining spatial features from
  video frames (VGG19/I3D-Kinetics) to preserve contextual details and temporally subsampling
  frames to improve efficiency. By fusing attended vectors across audio, video, and text modalities,
  the approach dynamically focuses on relevant cues during answer generation. This integrated
  framework addresses challenges in holistic dialog management, leveraging cross-modal interactions
  to outperform prior methods without relying on rigid pipelines, as demonstrated on the audio-visual
  scene-aware dataset.
- AVST [10]: The paper proposes a novel approach to Audio-Visual Question Answering (AVQA) by integrating multimodal understanding and spatio-temporal reasoning in dynamic audio-visual scenarios. It introduces the MUSIC-AVQA dataset with 45K QA pairs to benchmark the task, while addressing spatial associations through an attention-based sound source localization module (AV-Loc) to link sounds with visual sources. Temporal grounding (Q-Temp) is achieved by using question features to highlight key timestamps, enabling effective encoding of question-aware audio-visual embeddings. These components are fused to jointly represent spatial and temporal cues, overcoming challenges in cross-modal reasoning and enhancing performance in complex audio-visual scenes without relying on single-modality methods. The integrated framework demonstrates superior scene understanding by leveraging multisensory perception and fine-grained spatio-temporal analysis.
- AVSIAM [34]: The paper proposes an efficient and scalable audio-visual learning framework using a shared vision transformer backbone to unify audio and visual processing. The AVSiam model employs a contrastive audio-visual matching objective with a multi-ratio random masking scheme to enhance representation robustness while enabling larger batch sizes for effective contrastive learning. By sharing parameters across modalities, the approach reduces GPU memory footprint and computational costs compared to dual-backbone methods, while maintaining competitive performance on classification and retrieval tasks. This integrated design addresses scalability challenges and modality-handling flexibility without compromising accuracy.
- AMUSE [14]: The paper proposes a framework for music audio-visual question answering that
  addresses the unique challenges of dense, continuous audio-visual signals in musical performances.
  To exploit multimodal interconnectivity, it employs cross-modal adapters to facilitate early-stage
  token interactions between Swin-V2 (video), HTS-Audio (audio), and language transformers, while
  annotating rhythm and music sources in datasets to explicitly model musical characteristics. For
  temporal alignment, it designs specialized encoders to link musical signals with time dimensions.
  This integrated approach overcomes limitations of general-purpose AVQA methods by capturing intricate audio-visual relationships in performances, enhancing accuracy for music-specific questions
  through rhythm-aware and temporally grounded representations.

- ATT-BLSTM [85]: The paper proposes an attention-based bidirectional LSTM network (Att-BLSTM) for relation classification to capture decisive semantic information without relying on lexical resources or NLP systems. The model processes raw text through an embedding layer to generate word vectors, while bidirectional LSTM (BLSTM) layers learn high-level features by incorporating both past and future context. An attention mechanism then assigns weights to key words, merging word-level features into a sentence-level vector for classification. By integrating these components, the approach overcomes limitations of manual feature engineering and dependency on external tools, effectively identifying critical semantic cues across sentence positions to improve relation classification performance.
- AUDIO FLAMINGO [29]: The paper proposes Audio Flamingo, a novel audio language model designed to enhance large language models' (LLMs) understanding of non-speech sounds and non-verbal speech through three key innovations. It employs a sliding-window audio feature extractor to preserve temporal information in variable-length audio, while cross-attention mechanisms efficiently fuse audio inputs into the LM to reduce computational overhead. The model leverages a curated heterogeneous dataset and a two-stage training approach (pre-training and supervised fine-tuning) to balance close-ended and open-ended tasks. Additionally, it integrates in-context learning (ICL) and retrieval-augmented generation (RAG) through tailored templates and cross-attention masks, enabling few-shot adaptation without fine-tuning. To support multi-turn dialogues, the model is fine-tuned on GPT-4-generated datasets with correlated context. By combining these techniques, Audio Flamingo addresses challenges in audio feature extraction, heterogeneous data training, task adaptation, and dialogue coherence, achieving state-of-the-art performance across.
- CAT [13]: The paper proposes an enhanced Multimodal Large Language Model (MLLM) to improve question answering in dynamic audio-visual scenarios by addressing ambiguity and localization challenges. Key components include a clue aggregator to dynamically capture question-aware audio-visual features for fine-grained grounding, a mixed training strategy combining video-text and audio-text pairs with a novel AVinstruct dataset to strengthen cross-modal awareness, and an AI-assisted Ambiguity-aware Direct Preference Optimization (ADPO) to retrain the model for precise responses. By integrating these innovations, CAT effectively mitigates ambiguous outputs and enhances audio-visual reasoning, outperforming existing methods in Audio-Visual Question Answering (AVQA) tasks.
- CIGN [43]: The paper proposes a novel framework for continual audio-visual learning by disentangling class-aware cross-modal representations to mitigate catastrophic forgetting. It introduces learnable audio-visual class tokens to continually aggregate category-wise features through the Audio-Visual Continual Grouping module, while the Audio-Visual Class Tokens Distillation module preserves knowledge from previous tasks by aligning old and new token distributions. By integrating these components, the approach effectively addresses the challenge of mixed audio semantics and forgetting in sequential tasks, enhancing discriminative feature learning across modalities without relying on single-modality or rehearsal-based methods. The CIGN framework demonstrates superior performance in class-incremental audio-visual scenarios through its ability to maintain compact and disentangled representations.
- COCA [44]: The paper proposes a collaborative causal regularization framework (COCA) to address multi-shortcut biases in Audio-Visual Question Answering (AVQA) by integrating causal intervention and dynamic debiasing. The Bias-centered Causal Regularization (BCR) mitigates specific shortcut biases (Q→G, V&Q→G, A&Q→G) through counterfactual interventions to disrupt bias-irrelevant causal effects and factual regularization to maintain semantic consistency, while the Multi-shortcut Collaborative Debiasing (MCD) dynamically adjusts debiasing focus per sample using an entropy-driven metric to balance bias contributions. By jointly addressing unimodal and joint-modal biases through causal introspection and instance-aware adaptation, COCA enhances multimodal reasoning robustness without over-correcting, achieving state-of-the-art performance on MUSIC-AVQA.
- CONVLSTM [45]: The paper proposes a novel approach to enhance temporal reasoning in Audio Question Answering (AQA) by introducing the Diagnostic Audio Question Answering (DAQA) dataset, which comprises natural sound events and programmatically generated questions to probe temporal reasoning skills, while adapting visual question answering methods to AQA reveals their limitations. To address this, the authors develop Multiple Auxiliary Controllers for Linear Modulation (MALiMo), which extends Feature-wise Linear Modulation (FiLM) by incorporating an additional auxiliary controller to process subsampled audio features, thereby

enabling dynamic modulation of convolutional network processing based on both input modalities. This integrated approach improves relational and temporal reasoning by jointly leveraging audio and question inputs, overcoming the shortcomings of existing methods in handling complex temporal dependencies within sound sequences.

- CHATBRIDGE [39]: The paper proposes a multimodal language model that leverages large language models (LLMs) as a universal interface to bridge diverse modalities through language-paired data. ChatBridge integrates modality-specific encoders and perceiver modules to project embeddings into the LLM's semantic space, enabling cross-modal correlation without requiring all paired data combinations. The model undergoes two-stage training: first aligning modalities with language to emergent multimodal abilities, then instruction-finetuning on the MULTIS dataset to enhance zero-shot task generalization. By using language as a catalyst, ChatBridge addresses the challenge of limited multimodal paired data while achieving strong performance across text, image, video, and audio tasks through unified multimodal reasoning and user intent alignment.
- CROSSMAE [46]: The paper proposes a region-aware audio-visual pre-training framework to
  enhance cross-modality interaction and fine-grained alignment by extending masked autoencoders.
  It introduces Cross-Conditioned Reconstruction to reconstruct input pixels conditioned on crossmodal Attentive Tokens, while Cross-Embedding Reconstruction leverages Learnable Queries with
  positional cues to guide feature reconstruction between modalities, supplemented by contrastive
  loss for global alignment. By integrating these components, CrossMAE addresses the limitations
  of prior global feature-based methods, enabling effective region-level understanding and improving
  performance in both classification and dense prediction tasks without task-specific fine-tuning.
- DCL [49]: The paper proposes a disentangled counterfactual learning approach for physical audiovisual commonsense reasoning to infer objects' physics properties from video and audio inputs. The method decouples videos into static (time-invariant) and dynamic (time-varying) factors through a disentangled sequential encoder (DSE) using a variational autoencoder and contrastive loss to maximize mutual information while minimizing cross-factor interference. It further introduces a counterfactual learning module (CLM) to model physical knowledge relationships among objects by applying counterfactual interventions as confounders to enhance causal reasoning. By integrating DSE's disentangled representations with CLM's causal learning, the approach effectively addresses challenges in extracting implicit physical knowledge from multi-modal data, improving reasoning explainability and performance without relying on mixed feature representations.
- DG-SCT [52]: The paper proposes a novel Dual-Guided Spatial-Channel-Temporal (DG-SCT) attention mechanism to enhance large pre-trained models for audio-visual tasks by dynamically adjusting feature extraction through cross-modal guidance. The DG-SCT mechanism leverages audio and visual modalities as soft prompts to adaptively refine features across spatial, channel, and temporal dimensions, while preserving frozen pre-trained parameters. By integrating trainable cross-modal interaction layers into encoders, the approach emphasizes task-relevant information in each modality, addressing limitations of single-modality pre-training. This bidirectional prompting enables fine-grained feature fusion, improving performance on downstream tasks like AVE, AVVP, AVS, and AVQA without full retraining, while also excelling in few-shot and zero-shot scenarios.
- EEMC [53]: The paper proposes a novel task called Reference Audio-Visual Segmentation (Ref-AVS) to segment visual objects using expressions enriched with multimodal audio-visual cues, addressing the limitations of unimodal approaches. It introduces the Ref-AVS benchmark with pixellevel annotations and diverse multimodal-cue expressions to enable training and evaluation, while an end-to-end framework leverages a crossmodal transformer to process and integrate multimodal cues for precise segmentation. By simultaneously utilizing audio and visual descriptions in natural language, the approach overcomes challenges in locating objects in dynamic audio-visual scenes, enhancing segmentation accuracy in complex real-world scenarios without relying on manual mask annotations or single-modality references.
- FCNLSTM [45]: The paper proposes a novel approach to enhance temporal reasoning in Audio Question Answering (AQA) by introducing the Diagnostic Audio Question Answering (DAQA) dataset, which comprises natural sound events and programmatically generated questions to probe temporal reasoning skills. While adapting existing visual question answering methods to AQA reveals their limitations in temporal reasoning, the authors develop Multiple Auxiliary Controllers for Linear Modulation (MALiMo) to extend Feature-wise Linear Modulation (FiLM) by incorporating an additional auxiliary controller to process subsampled audio features, thereby enabling dynamic modulation of convolutional network processing based on both principal and

supplementary inputs. This integrated approach addresses the challenge of in-depth temporal reasoning by facilitating relational and temporal analysis, leading to improved performance on DAQA without relying on spatial reasoning or static inputs.

- GPT-40 [4]: The paper proposes GPT-40, an autoregressive omni model designed to process any combination of text, audio, image, and video inputs while generating text, audio, or image outputs through end-to-end training across modalities. By integrating Web Data, Code and Math, and Multimodal Data during pre-training, the model learns diverse reasoning skills and multimodal interpretation, while post-training alignment and red-teaming mitigate risks such as bias and harmful content. This unified approach enhances real-time responsiveness, multilingual performance, and multimodal understanding while addressing safety concerns through layered mitigations and external evaluations.
- GRU [55]: The paper proposes a free-form, open-ended Visual Question Answering (VQA) task
  to generate natural language answers from images and questions, mirroring real-world scenarios
  like assisting the visually impaired. The approach leverages a large dataset (0.25M images, 0.76M
  questions, 10M answers) combining real images from MS COCO and abstract scenes to enable both
  low-level vision and high-level reasoning. By supporting diverse question types (e.g., fine-grained
  recognition, commonsense reasoning) and offering automatic evaluation through open-ended or
  multiple-choice formats, the framework addresses the need for detailed image understanding and
  multi-modal knowledge integration, advancing AI-complete challenges beyond generic captioning.
- HCATTN [86]: The paper proposes a hierarchical co-attention model for Visual Question Answering (VQA) that jointly reasons about image and question attention to improve answer accuracy. It introduces a co-attention mechanism to simultaneously perform question-guided visual attention (to identify relevant image regions) and image-guided question attention (to focus on key words), while employing a hierarchical question representation through word-level embeddings, phrase-level 1D CNNs (to capture n-gram features), and question-level LSTMs (to encode contextual meaning). By recursively combining co-attended features across these levels, the model addresses challenges like linguistic variation and multi-modal alignment, enhancing robustness and fine-grained understanding for VQA tasks.
- HCRN [57]: The paper proposes a general-purpose neural unit for video question answering that enables hierarchical relational reasoning and multimodal fusion. The Conditional Relation Network (CRN) processes input object arrays through sparse high-order relations while modulating encodings with contextual features, allowing flexible replication and stacking into Hierarchical CRNs (HCRN). The architecture integrates appearance features with clip motion as initial context, then progressively incorporates linguistic context and video-level motion through layered CRNs to enable multi-step reasoning. By hierarchically combining localized clip relations with global video and question contexts, HCRN addresses challenges of modeling distant temporal dependencies and heterogeneous modalities in VideoQA, demonstrating robust performance across diverse question types requiring appearance, motion, and temporal reasoning.
- HME [87]: The paper proposes a novel VideoQA framework that integrates heterogeneous memory and multimodal attention to enhance video-question reasoning. It introduces a heterogeneous memory module to jointly learn global context from appearance and motion features through synchronized attention, while a redesigned question memory captures complex semantics and highlights queried subjects by storing global contexts. These components interact through a multimodal fusion layer that aligns visual and textual hints via self-updated attention, enabling multi-step reasoning. By unifying feature integration with attention learning and maintaining global context throughout, the approach addresses challenges of spatiotemporal alignment and complex question semantics, improving VideoQA performance without separating feature and attention steps.
- LAST-ATT [11]: The paper proposes a method to address data bias in audio-visual question answering (AVQA) by constructing a balanced dataset and introducing an enhanced multimodal model. It identifies skewed answer distributions in the MUSIC-AVQA dataset and rectifies them by collecting complementary videos and questions to ensure uniform answer spread, particularly for binary questions, resulting in the MUSIC-AVQA v2.0 benchmark. The baseline model strengthens audio-visual-text interrelations through a pretrained Audio-Spectrogram-Transformer (AST) branch for audio grounding and cross-modal pixel-wise attention to align audio and visual spatial maps. By integrating these components, the approach mitigates modality neglect and improves reasoning across vision, audio, and language, establishing a robust foundation for unbiased AVQA evaluation.

- LAVIT [59]: The paper proposes a novel benchmark for grounded audio-visual question answering on 360° videos to address spherical spatial reasoning and audio-visual relationships. It introduces the Pano-AVQA dataset with 51.7K QA pairs, featuring bounding-box grounding for two task types: spherical spatial relation QAs to assess relative object positioning on a sphere, and audio-visual relation QAs to link sounds with visual sources. Through quaternion-based spatial embeddings and multimodal training objectives, the framework integrates panoramic audio-visual cues while addressing challenges like spherical distortion and diverse sound localization. This holistic approach enhances semantic understanding of omnidirectional environments without relying on predefined fields of view.
- LAVISH [58]: The paper proposes adapting frozen vision transformers (ViTs) pretrained on visual data to audio-visual tasks without finetuning their original parameters. This is achieved through a latent audio-visual hybrid (LAVISH) adapter, which injects trainable parameters into each ViT layer to enable audio specialization and cross-modal fusion. The LAVISH adapter employs latent tokens to compress modality-specific information, reducing the quadratic cost of standard cross-attention while facilitating bidirectional audio-visual interaction. By integrating these components, the approach addresses the inefficiency of modality-specific models and costly audio pretraining, enabling frozen ViTs to leverage shared representations for enhanced audio-visual understanding without external encoders or extensive parameter updates.
- LSTTA [60]: The paper proposes a parameter-efficient transfer learning approach for audiovisual-language tasks by introducing the Long Short-Term Trimodal Adapter (LSTTA), which integrates pre-trained unimodal/bimodal models without full fine-tuning. LSTTA employs a longterm semantic filtering module to suppress redundant video frames by characterizing temporal importance, while the short-term semantic interaction module models local cross-modal alignments through two sub-modules (AL2V and VL2A) to facilitate fine-grained information transfer. By combining these complementary mechanisms, LSTTA addresses the challenges of uneven global semantics and unannotated local correspondences in trimodal learning, enhancing performance on tasks like Music-AVQA and CMU-MOSEI without requiring large-scale trimodal pretraining.
- MAVEN [62]: The paper proposes a robust multimodal reasoning framework for Audio-Visual Question Answering (AVQA) to address dataset biases and enhance model robustness. It introduces FortisAVQA, a novel dataset constructed by rephrasing test questions to diversify linguistic forms and introducing distribution shifts to evaluate performance across frequent and rare question types. The Multimodal Audio-Visual Epistemic Network (MAVEN) employs a Multifaceted Cycle Collaborative Debias (MCCD) strategy to mitigate bias learning by enlarging distribution differences between unimodal and multimodal logits through KL divergence optimization while using cycle guidance to align unimodal logit distributions. This integrated approach reduces reliance on spurious correlations in individual modalities, improving generalization across in-distribution and out-of-distribution scenarios without requiring balanced training data.
- MCAN [64]: The paper proposes a deep Modular Co-Attention Network (MCAN) to enhance visual question answering (VQA) by jointly modeling intra- and inter-modal interactions through a modular architecture. The framework integrates Self-Attention (SA) units to capture dense word-to-word and region-to-region relationships within questions and images, while Guided-Attention (GA) units model word-to-region cross-modal dependencies. By cascading Modular Co-Attention (MCA) layers composed of SA and GA units, MCAN enables deep reasoning while addressing the limitations of shallow co-attention models. This integrated approach improves fine-grained semantic understanding by simultaneously refining self-attention within modalities and guided-attention across modalities, leading to more accurate visual-textual alignment and robust performance on complex VQA tasks.
- MCCD [23]: The paper proposes a robust framework for Audio-Visual Question Answering (AVQA) to address dataset biases and enhance model robustness. It introduces MUSIC-AVQA-R, a novel dataset crafted by rephrasing test questions and introducing distribution shifts to evaluate performance on both frequent and rare samples, while the Multifaceted Cycle Collaborative Debiasing (MCCD) strategy mitigates bias learning by enlarging distribution differences between uni-modal and multi-modal logits and employing cycle guidance to align uni-modal distributions. This integrated approach ensures diverse question evaluation and reduces bias dependency, improving generalization across in- and out-of-distribution scenarios without relying on balanced training data.

- MEERKAT [67]: The paper proposes an audio-visual LLM for fine-grained spatio-temporal grounding in images and audio, addressing the limitations of existing MLLMs in handling fine-grained tasks. It introduces a modality alignment module based on optimal transport to learn cross-modal patch alignment in a weakly-supervised manner, while a cross-attention module enforces audiovisual consistency to improve joint representation learning. These components are integrated through the AVFIT dataset (3M instruction samples) and MeerkatBench, a unified benchmark for five tasks, enabling the model to tackle challenges like disparate task formats and lack of large-scale training data. The approach enhances performance by unifying spatial and temporal grounding capabilities, achieving state-of-the-art results across diverse audio-visual tasks.
- OPM [69]: The paper proposes an adaptive modulation approach to address imbalanced multimodal learning by dynamically balancing uni-modal optimization during joint training. It introduces Onthe-fly Prediction Modulation (OPM) to weaken dominant modality influence in the feed-forward stage by probabilistically dropping its features, while On-the-fly Gradient Modulation (OGM) mitigates gradient dominance in back-propagation through adaptive noise injection. By monitoring inter-modal discriminative discrepancies, these strategies jointly alleviate under-optimization of weaker modalities while preserving dominant modality contributions. The integrated framework enhances multimodal representation learning across diverse tasks by ensuring balanced feature optimization without additional training overhead, as validated through extensive experiments on audio-visual benchmarks.
- ONELLM [2]: The paper proposes a unified framework to align multiple modalities with language using a shared architecture, eliminating the need for modality-specific encoders. It introduces lightweight modality tokenizers to convert input signals into tokens, while a universal encoder (CLIP-ViT) extracts cross-modal features and a universal projection module (UPM) dynamically routes mixed projection experts to map diverse modalities into the LLM's embedding space. Through progressive alignment and a curated multimodal instruction dataset spanning eight modalities, the integrated approach overcomes scalability limitations of prior MLLMs by unifying encoding and projection, enabling flexible modality expansion and enhanced multimodal understanding without architectural redundancy.
- PSAC [70]: The paper proposes a novel self-attention-based architecture for video question answering (VQA) to overcome the limitations of RNNs in modeling long-range dependencies and parallel processing. It introduces Positional Self-Attention (PSA) to capture global dependencies in video and question sequences by attending to all positions while incorporating absolute positional encodings to preserve temporal/spatial information. Through Video-based PSA (VPSA) and Question-based PSA (QPSA), the model encodes video frames and textual questions in parallel. A Video-Question Co-Attention (VQ-Co) block then simultaneously attends to relevant visual and textual features via bidirectional attention, enhancing cross-modal alignment. By integrating PSA with co-attention, the framework efficiently models complex video-question interactions, addressing challenges in sequential data processing and multimodal fusion while improving accuracy and computational efficiency.
- PSTP-NET [25]: The paper proposes a progressive spatio-temporal perception framework for audio-visual question answering (AVQA) to address challenges in complex multi-modal video understanding. The Temporal Segment Selection Module (TSSM) identifies relevant video segments to reduce redundancy, while the Spatial Region Selection Module (SRSM) locates question-aware visual patches within selected segments to enhance spatial reasoning. The Audio-guided Visual Attention Module (AVAM) models audio-visual associations by aligning sound features with visual patches. By progressively integrating these components, the approach effectively filters irrelevant content, localizes key spatio-temporal regions, and strengthens cross-modal interactions, leading to improved scene understanding and question answering performance.
- QAP [5]: The paper proposes a parameter-efficient multimodal language model learning strategy that bridges modalities through query-based prompts and lightweight resampling. The core innovation involves Querying Prompts (QP) to simultaneously extract modality information and interact with text, while Text-Conditioned Resamplers (TCR) adaptively inject text-relevant multimodal features into frozen language model layers. By integrating QP and TCR, the approach efficiently compresses modality inputs and leverages the model's inherent fusion capabilities, addressing computational inefficiency and redundancy in traditional projection-based methods while outperforming existing techniques across multiple multimodal tasks with minimal trainable parameters.

- QWEN2.5-VL [3]: The paper proposes Qwen2.5-VL, a vision-language model advancing multimodal understanding through enhanced visual recognition, object localization, and document parsing while addressing computational and contextual challenges. Key innovations include dynamic resolution processing to handle varying image sizes and video durations, absolute time encoding to improve temporal dynamics perception, and a native dynamic-resolution ViT with Window Attention to reduce overhead while preserving resolution. By integrating these components, the model achieves robust performance in fine-grained visual tasks, long-video comprehension, and real-world agentic applications without task-specific fine-tuning, while maintaining strong linguistic capabilities inherited from Qwen2.5 LLM. The approach overcomes bottlenecks in computational complexity and inconsistent sequence-length performance, enabling precise spatial-temporal reasoning and cross-domain generalization.
- REFATOMNET [71]: The paper proposes a novel approach for Referring Atomic Video Action Recognition (RAVAR) to identify atomic actions of a specific person guided by textual descriptions and video data, addressing limitations in traditional action recognition. Key components include RefAtomNet, which employs a multi-stream architecture connecting video, text, and location-semantic streams to interpret referring expressions and localize target individuals, while cross-stream agent attention and token fusion enhance relevance filtering across modalities. This integrated approach tackles challenges like irrelevant visual distractions and enables end-to-end action recognition for referred individuals, outperforming existing methods in RAVAR without requiring manual pre-processing. The RefAVA dataset with 36,630 annotated instances supports this task.
- VALOR [72]: The paper proposes a Vision-Audio-Language Omni-Perception pretraining model (VALOR) to jointly model tri-modality interactions for understanding and generation tasks. It employs three single-modality encoders to process vision, audio, and language separately, while a multimodal decoder enables conditional text generation through two pretext tasks: Multimodal Grouping Alignment (MGA) projects modalities into a shared space to align vision-language, audiolanguage, and audiovisual-language groups via contrastive learning, and Multimodal Grouping Captioning (MGC) reconstructs masked text tokens conditioned on visual, auditory, or combined inputs to enhance generative capabilities. By integrating these components with a large-scale human-annotated dataset (VALOR-1M), the approach addresses the limitations of existing bimodal systems, enabling comprehensive cross-modal alignment and flexible text generation across diverse modality combinations for downstream tasks like retrieval, captioning, and question answering.
- VAST [73]: The paper proposes an omni-modality foundation model to enhance video-text crossmodality learning by integrating vision, audio, and subtitle information. It introduces VAST-27M, a large-scale dataset automatically generated through a two-stage pipeline: first training separate vision and audio captioners to produce single-modality descriptions, then employing an LLM to synthesize these with subtitles into omni-modality captions. The VAST model leverages three modality encoders and cross-attention-based text fusion, trained with objectives (OM-VCC/VCM/VCG) to unify multi-modal understanding. This approach addresses the lack of comprehensive video-text corpora by automating caption generation, enabling joint modeling of complementary modalities to improve performance on diverse downstream tasks like retrieval, captioning, and QA without manual annotation costs.
- VITA [77]: The paper proposes VITA, an open-source Multimodal Large Language Model (MLLM) capable of simultaneous processing and interactive analysis across video, image, text, and audio modalities. Starting with Mixtral 8×7B as a language foundation, it expands Chinese vocabulary through bilingual instruction tuning to enhance multilingual proficiency, while endowing visual and audio capabilities via two-stage multi-task learning for multimodal alignment and instruction tuning. To improve interaction, VITA introduces state tokens to distinguish input queries for non-awakening interaction and employs a duplex pipeline deployment scheme, where one model generates responses while another monitors environmental inputs, enabling audio interrupt interaction. This integrated approach addresses the lack of open-source models with unified multimodal processing and natural interaction, advancing seamless multimodal understanding and human-computer engagement without relying on wake-up words or sequential query handling.
- VIDEOLLAMA-2 [76]: The paper proposes VideoLLaMA 2, a Video Large Language Model designed to enhance spatial-temporal modeling and audio understanding in multimodal video tasks. It introduces a Spatial-Temporal Convolution (STC) connector to capture intricate spatial and temporal dynamics in video data, while integrating an Audio Branch through joint training to

incorporate audio cues for richer multimodal understanding. By combining these components, the model addresses challenges in processing temporal dynamics and audio-visual synchronization, improving performance in video question answering and captioning tasks without compromising contextual integrity or processing efficiency.