

# Distorted or Fabricated? A Survey on Hallucination in Video LLMs

Anonymous ACL submission

## Abstract

Despite significant progress in video-language modeling, hallucinations remain a persistent challenge in Video Large Language Models (Vid-LLMs), referring to outputs that appear plausible yet contradict the content of the input video. This survey presents a comprehensive analysis of hallucinations in Vid-LLMs and introduces a systematic taxonomy that categorizes them into two core types: dynamic distortion and content fabrication, each comprising two subtypes with representative cases. Building on this taxonomy, we review recent advances in the evaluation and mitigation of hallucinations, covering key benchmarks, metrics, and intervention strategies. We further analyze the root causes of dynamic distortion and content fabrication, which often result from limited capacity for temporal representation and insufficient visual grounding. These insights inform several promising directions for future work, including the development of motion-aware visual encoders and the integration of counterfactual learning techniques. This survey consolidates scattered progress to foster a systematic understanding of hallucinations in Vid-LLMs, laying the groundwork for building robust and reliable video-language systems.

## 1 Introduction

Video Large Language Models (Vid-LLMs) extend the capabilities of vision-language systems from static images to temporally coherent video inputs, enabling tasks such as action recognition, temporal reasoning, and audio-visual understanding (Zhang et al., 2023a; Maaz et al., 2024; Li et al., 2023a; Lin et al., 2024; Wang et al., 2024a; Fu et al., 2024). Despite recent advances, these models remain susceptible to hallucinations, producing outputs that appear plausible and coherent yet contradict the actual content of the video. This issue poses reliability and safety risks in safety-critical domains, including embodied AI (Wu et al., 2023) and autonomous driving (Chen et al., 2024).

While hallucinations have been extensively surveyed in image-based vision-language models (VLMs) (Liu et al., 2024a; Lan et al., 2024), the inherent complexity of video’s temporal structure, motion dynamics, and audio-visual integration complicates the direct application of these insights to the video domain. To address this gap, this survey presents a video-specific, mechanism-driven taxonomy that classifies hallucinations into two primary types: dynamic distortion, where the model misrepresents the spatiotemporal evolution or referential consistency of entities and scenes; and content fabrication, where outputs are influenced by prior knowledge or dominated by audio modality.

Building on this taxonomy, we review recent advances in the evaluation and mitigation of hallucinations in Vid-LLMs, with a focus on key benchmarks, metrics, and intervention strategies. Dynamic distortion includes hallucinations in spatiotemporal dynamics, such as incorrect event ordering (Li et al., 2025a; Wu et al., 2025a; Sun et al., 2025), inaccurate duration estimation (Wang et al., 2024b; Huang et al., 2025d; Sun et al., 2024), and frequency miscounting (Gao et al., 2025; Choong et al., 2024), as well as referential inconsistency, where the model conflates different characters (Seth et al., 2025; Yang et al., 2025) or scenes (Lu et al., 2025; Ma et al., 2024; Pu et al., 2025). Content fabrication includes context-driven hallucinations, where commonly co-occurring object–action (Chang et al., 2025; Li et al., 2025c) or scene–event (Bae et al., 2025; Zhang et al., 2024; Ding et al., 2025) patterns lead to unsupported inferences; and audio-visual conflict, where dominant auditory cues override visual evidence, resulting in hallucinated actions (Sung-Bin et al., 2025; Jung et al., 2025) or emotional states (Xing et al., 2025).

Further analysis reveals the underlying mechanisms of dynamic distortion and content fabrication. Dynamic distortion often results from missing fine-grained motion cues due to limited temporal

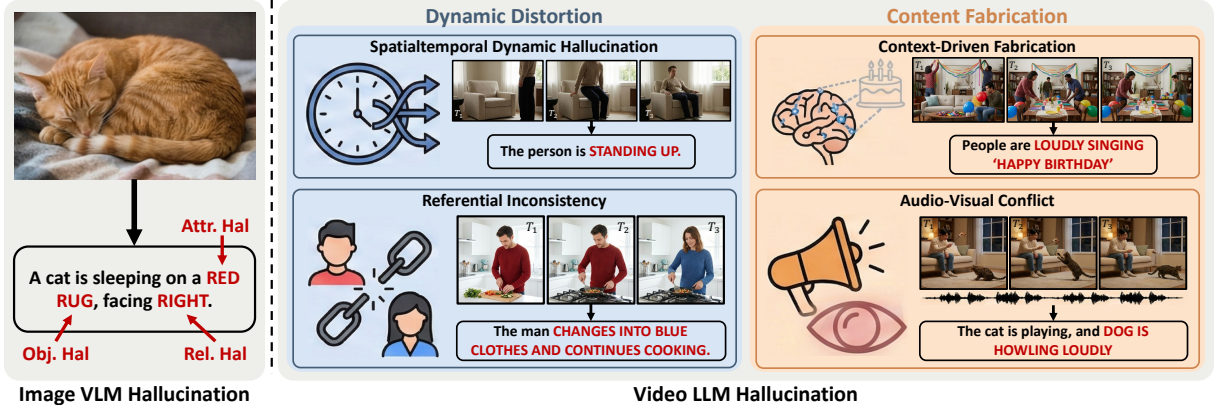


Figure 1: Taxonomy of hallucinations in Vid-LLMs. Video models exhibit unique hallucination types beyond static errors in images. These fall into two main categories: (1) Dynamic Distortion, which includes spatiotemporal misrepresentation and referential inconsistency; (2) Content Fabrication, which includes hallucinations influenced by statistical priors and cases where auditory input overrides visual evidence.

encoding (Zhao et al., 2025; Liu et al., 2024b), and is further exacerbated in long videos by weak long-range memory (Bae et al., 2025) and poor temporal localization (Wu et al., 2025a). In contrast, content fabrication arises from insufficient visual grounding (Lee et al., 2025), allowing pretrained priors (Li et al., 2025c) or dominant audio signals (Leng et al., 2024) to override visual evidence.

In light of these underlying mechanisms, promising research directions include developing motion-aware architectures (Wu et al., 2025b) that retain fine-grained temporal features to strengthen the alignment between visual perception and temporal reasoning. In addition, counterfactual training strategies that disentangle visual evidence from prior knowledge (Huang et al., 2025c) offer a principled approach to mitigating content fabrication by encouraging models to ground predictions more faithfully in the visual input.

**Comparison with existing surveys.** Hallucination has been extensively studied in LLMs and image-based VLMs (Zhang et al., 2023c; Huang et al., 2025b; Liu et al., 2024a; Lan et al., 2024). While MLLM hallucination surveys (Sahoo et al., 2024; Bai et al., 2024) include video alongside other modalities, their discussion of video hallucination remains superficial, offering only brief mentions of benchmarks and mitigation strategies without structural or causal analysis. In contrast, this survey presents the first mechanism-driven taxonomy of hallucinations in Vid-LLMs. We propose a layered classification framework (Fig. 2), conduct a broader and more detailed review of existing literature, and analyze the underlying causes of hallucina-

tions. Building on this analysis, we outline future directions that align closely with identified causes, benchmark coverage, and mitigation strategies, offering a cohesive roadmap toward hallucination-resilient Vid-LLMs.

## 2 Definition and Scope

**Definition.** We define *video hallucination* as cases where a Vid-LLM generates textual outputs that are linguistically coherent and contextually plausible, yet contradict the observable spatiotemporal evidence in the input video.

**Distinction from Static Image Hallucination.** While hallucinations in image-based VLMs, such as those involving objects, attributes, and relations, have been well studied (Liu et al., 2024a; Lan et al., 2024), the video modality introduces a temporal dimension that fundamentally alters the problem. Unlike static inputs, videos require reasoning over causality, temporal grounding, motion dynamics, and audio-visual integration. These aspects are beyond the scope of traditional static metrics. Our taxonomy explicitly captures these temporal and multimodal challenges, distinguishing video hallucination from image-based settings.

## 3 Taxonomy of Video Hallucinations

As discussed in Section 2, the temporal and multimodal nature of video poses challenges beyond static image settings. To address these, we propose a mechanism-driven taxonomy focused on *dynamic-level hallucinations* unique to video. It captures how Vid-LLMs fail in temporal reasoning and cross-modal alignment. We classify hallucina-

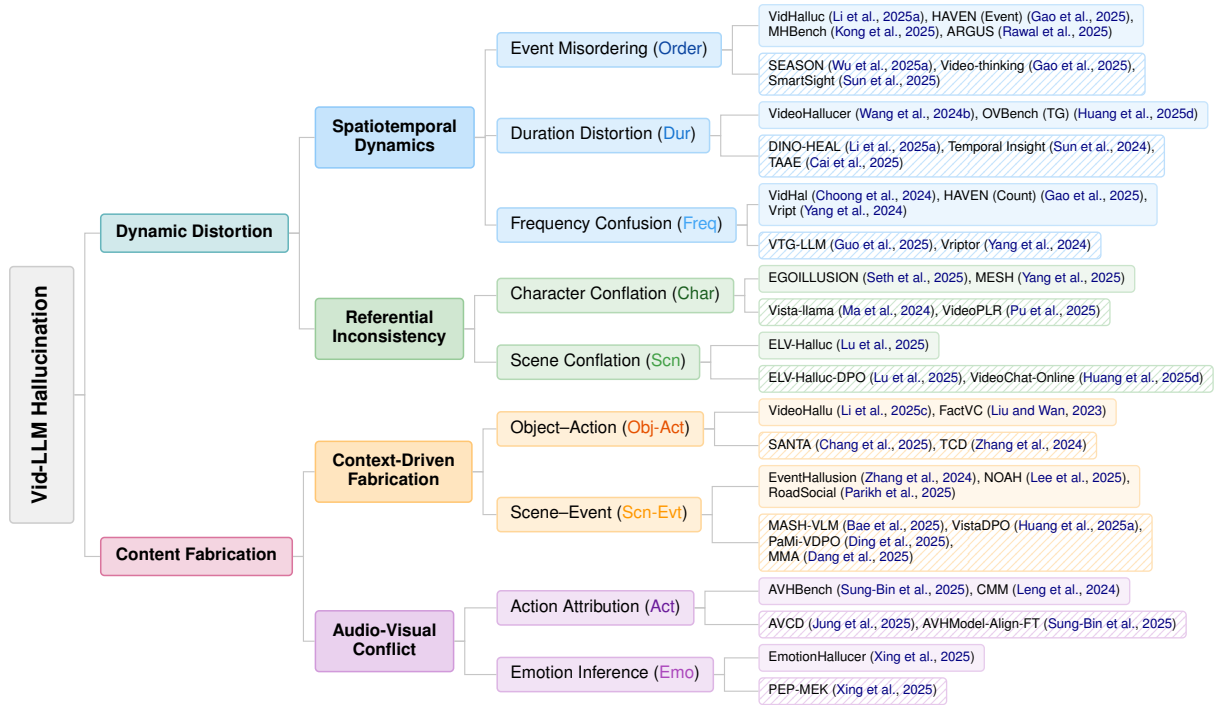


Figure 2: Mechanism-driven taxonomy of Vid-LLM hallucinations. **Dynamic Distortion**: entities are perceived but their spatiotemporal evolution or identity is misinterpreted, including **Spatiotemporal Dynamics** (Order/Dur/Freq) and **Referential Inconsistency** (Char/Scn). **Content Fabrication**: outputs lack visual evidence and are driven by priors, including **Context-Driven Fabrication** (Obj-Act/Scn-Evt) and **Audio-Visual Conflict** (Act/Emo). Solid fill denotes benchmarks; striped fill indicates mitigation methods.

tions into *Dynamic Distortion* and *Content Fabrication*, each with two subtypes and representative cases (Figure 2). This taxonomy provides a foundation for structuring benchmarks (Section 4) and mitigation strategies (Section 5).

### 3.1 Dynamic Distortion

This category refers to situations where the model correctly detects entities but misrepresents their temporal progression or referential consistency. It includes two subtypes: *Spatiotemporal Dynamics*, involving errors in event ordering, duration, or frequency; and *Referential Inconsistency*, where characters or scenes are conflated across temporal boundaries.

**Spatiotemporal Dynamics.** These hallucinations arise when the model correctly identifies relevant events but fails to model their temporal relationships. Typical cases include event misordering, such as reversing action causality or misinterpreting motion direction and trajectory (Li et al., 2025a; Gao et al., 2025; Wu et al., 2025a; Sun et al., 2025); duration distortion, where the model over- or underestimates the length of an action (Wang et al., 2024b; Huang et al., 2025d; Sun et al., 2024); and frequency confusion, in which repeated actions are

miscounted (Gao et al., 2025; Choong et al., 2024).

**Referential Inconsistency** These hallucinations refers to semantic-level failures where the model conflates distinct entities or scenes across temporal boundaries, producing blended descriptions that obscure segment distinctions. These errors arise when content from separate time spans is incorrectly merged into a single entity- or scene-level statement, even when visual cues could distinguish them. Such inconsistency typically appears in two forms: character conflation, where different individuals across scenes are mistakenly treated as the same person (Seth et al., 2025; Yang et al., 2025); and scene conflation, where actions or settings from distinct contexts are combined into a single narrative (Lu et al., 2025; Pu et al., 2025).

### 3.2 Content Fabrication

This category covers cases where the model produces outputs that lack grounding in visual evidence and are instead influenced by learned priors. It includes *context-driven fabrication*, where common object-action or scene-event associations result in unsupported predictions, and *audio-visual conflict*, where auditory cues override visual input.

**Context-Driven Fabrication** This type of hallu-

cination arises when the model relies on statistical associations from training data rather than grounding its predictions in visual evidence. An error is considered context-driven fabrication when the predicted action or event lacks visual support in the current observation window but is triggered by the presence of associated objects or scenes. It typically appears in two forms: object–action fabrication and scene–event fabrication. Object–action fabrication (Chang et al., 2025; Li et al., 2025c; Liu and Wan, 2023) occurs when the presence of an object leads to incorrect action inference despite lacking motion cues. Scene–event fabrication (Bae et al., 2025; Zhang et al., 2024; Ding et al., 2025; Dang et al., 2025; Huang et al., 2025a) happens when typical events are predicted solely from background settings.

**Audio-Visual Conflict** This type of fabrication occurs when dominant or misleading audio cues override visual evidence, leading the model to generate outputs that align more with the audio than the video. Typical cases include hallucinated actions triggered by background sounds (Sung-Bin et al., 2025; Jung et al., 2025), and emotion inference based on vocal tone rather than facial expression (Xing et al., 2025).

## 4 Evaluation Benchmarks

Following the taxonomy in Section 3, we categorize existing benchmarks by hallucination type and representative failure cases. Table 1 provides an overview of their venues, scales, task formats, and evaluation metrics.

### 4.1 Dynamic Distortion Benchmarks

**Spatiotemporal dynamics benchmarks** assess Vid-LLMs’ ability to model temporal structure, covering three subtypes: event misordering, duration distortion, and frequency confusion.

For *event misordering*, VidHalluc (Li et al., 2025a) includes 5,002 videos from ActivityNet (Yu et al., 2019), YouCook2 (Zhou et al., 2018), and VALOR32K (Liu et al., 2025a), and evaluates temporal hallucinations through sequence-based QA tasks that test whether models can determine the correct order of actions. HAVEN (Gao et al., 2025) (event) targets discrepancies in action sequences using 2,245 questions across binary, multiple-choice, and short-answer formats. MHBench (Kong et al., 2025) provides 1,200 videos and tests motion understanding via adversarial triplets simulating

original, reversed, and incomplete actions. AR-GUS (Rawal et al., 2025) evaluates hallucination and omission on 500 videos with about 9,500 annotations, penalizing event misordering by checking the temporal alignment between model-generated and ground-truth action sequences.

For *duration distortion*, VideoHalluc (Wang et al., 2024b) includes 1,800 adversarial question pairs based on 948 videos, assessing both intrinsic and extrinsic hallucinations through tasks focused on detecting abnormal durations and comparing relative event lengths. OVBench (THV) (Huang et al., 2025d) targets duration distortion in real-time streaming settings, requiring models to track action persistence and estimate the length of ongoing events as temporal context unfolds.

For *frequency confusion*, VidHal (Choong et al., 2024) benchmarks fine-grained temporal understanding by asking models to distinguish between captions with correct and hallucinated action counts. HAVEN (Gao et al., 2025) (count) addresses this via numerical questions that test a model’s ability to differentiate between single and repeated actions, evaluating its capacity to quantify frequency. Vript (Yang et al., 2024) includes a ‘Count’ category in its Vript-RR benchmark, which assesses whether models can accurately compare the number of visual elements across long video sequences.

**Referential inconsistency benchmarks** assess whether models can maintain distinct representations of entities and scenes over time. Despite the increasing use of Vid-LLMs, only three benchmarks explicitly address this issue.

For *character conflation*, EGOILLUSION (Seth et al., 2025) includes 1,400 egocentric videos and 8,000 question–answer pairs. It evaluates whether models confuse different individuals, for example by identifying the camera wearer as another person during object interactions or activity recognition. MESH (Yang et al., 2025) introduces a human-aligned evaluation framework called Mise En Scène, built on TVQA+ clips. It tests whether models can consistently track character identity, appearance, and actions across scenes using structured evaluation traps.

For *scene conflation*, ELV-Halluc (Lu et al., 2025) uses approximately 8,600 adversarial video and text pairs to evaluate whether models incorrectly assign visual elements such as objects or actions from one part of a video to another.

Table 1: Summary of video hallucination benchmarks. *Format*: **MC** = Multiple Choice, **Bin** = Yes/No, **Open** = Open-ended QA, **Cap** = Captioning. *Length*: **S** = Short (<1min), **M** = Medium (1–5min), **L** = Long (>5min), **St** = Streaming. *Baseline*: Specialized baseline method proposed. *SOTA Perf.*: Representative best performance reported.

Benchmark	Venue	# Vid	# QA	Format	Metric	Len	Domain	Baseline	SOTA Perf.
<i>Spatiotemporal Dynamics Benchmarks (Dynamic Distortion)</i>									
VidHalluc (Li et al., 2025a)	CVPR'25	5,002	9,295	<b>MC, Bin</b>	Acc, Score	<b>S</b>	ActivityNet, YouCook2, VALOR	✓	GPT-4o: 81.2%
VideoHalluc (Wang et al., 2024b)	arXiv'24	948	1,800	<b>Bin</b>	Acc, Score	<b>M</b>	ActivityNet, VidOR, YouCook	✓	Gemini-1.5: 37.8%
HAVEN (Gao et al., 2025)	arXiv'25	–	6.5k	<b>MC, Bin, Open</b>	Acc, Bias	<b>S</b>	COIN, ActivityNet, Sports1M	✓	Valley-Eagle: 61.3%
MHBench (Kong et al., 2025)	AAAI'25	1,200	–	<b>MC, Bin</b>	Acc, F1	<b>S</b>	Sth-Sth V2, Self-shot	✓	VideoChat2-MCD: 65.2%
VidHal (Choong et al., 2024)	arXiv'24	1,000	3,000	<b>MC</b>	Acc, NDCG	<b>S</b>	TempCompass, MVBench, PT	✗	GPT-4o: 77.2%
ARGUS (Rawal et al., 2025)	arXiv'25	500	~9.5k	<b>Cap</b>	Cost-H/O	<b>S</b>	Ego4D, Panda-70M, Stock	✗	Gemini-2.0: 41%
OVBench (THV) (Huang et al., 2025d)	CVPR'25	–	~33k	<b>Bin</b>	Accuracy	<b>St</b>	DiDeMo, QuerYD.	✓	VideoChat-On: 63.1%
Vript (Yang et al., 2024)	NeurIPS'24	12k	420k	<b>Bin, MC</b>	Acc, F1	<b>L</b>	HD-VILA, YouTube, TikTok	✓	Vriptor: 58.3 (F1)
<i>Referential Inconsistency Benchmarks (Dynamic Distortion)</i>									
EGOILLUSION (Seth et al., 2025)	EMNLP'25	1,400	8k	<b>Bin, Open</b>	Accuracy	<b>S/M</b>	Ego4D, EgoSeg, Trek-150	✗	Gemini-Pro: 59.4%
MESH (Yang et al., 2025)	MM'25	–	~140k	<b>MC, Bin</b>	Accuracy	<b>S/M</b>	TVQA+, UCF101	✗	GPT-4o: 79.1%
ELV-Halluc (Lu et al., 2025)	arXiv'25	200	4.8k	<b>Bin</b>	Acc, SAH	<b>L</b>	YouTube (Event-based)	✓	Gemini2.5-Flash: 53.1%
<i>Context-Driven Fabrication Benchmarks (Content Fabrication)</i>									
FactVC (Liu and Wan, 2023)	EMNLP'23	300	–	<b>Cap</b>	Bleu4, Rouge-L	<b>M/L</b>	ActivityNet, YouCook2	✓	PDVC-gt: 12.83 (Bleu4)
EventHallusion (Zhang et al., 2024)	AAAI'26	400	711	<b>Bin, Open</b>	Accuracy	<b>S</b>	ActivityNet	✓	GPT-4o: 91.93%
NOAH (Lee et al., 2025)	arXiv'25	9k	~60k	<b>Bin, Cap</b>	Acc, HR	<b>M/L</b>	ActivityNet	✗	Gemini2.5-Flash: 66.8%
VideoHallu (Li et al., 2025c)	NeurIPS'25	987	3,233	<b>Open</b>	GPT-Score	<b>S</b>	Generated (Sora, etc.)	✓	Comb-GRPO: 57.7
RoadSocial (Parikh et al., 2025)	CVPR'25	13.2k	260k	<b>Open</b>	GPT-Score	<b>S/M</b>	Social Media (Traffic)	✗	GPT-4o: 69.8
<i>Audio-Visual Conflict Benchmarks (Content Fabrication)</i>									
AVHBench (Sung-Bin et al., 2025)	ICLR'25	2,136	5.3k	<b>Bin</b>	Accuracy	<b>S</b>	AudioCaps, VALOR	✓	AVHModel-Align-FT: 83.9%
CMM (Leng et al., 2024)	arXiv'24	1.2k	2.4k	<b>Bin</b>	PA/HR	<b>S</b>	WebVid, AudioCaps	✗	Gemini-1.5: 88.4/64.2
EmotionHalluc (Xing et al., 2025)	arXiv'25	230	2,742	<b>Bin</b>	Accuracy	<b>S/M</b>	MER 2023, Social-IQ 2.0	✓	Gemini2.5-Flash: 68.2%

## 4.2 Content Fabrication Benchmarks

**Context-driven fabrication benchmarks** assess whether models generate outputs based on visual evidence rather than relying on statistical associations from training data. These benchmarks span various domains such as activity recognition, driving, and synthetic videos, reflecting the diverse and context-sensitive nature of fabrication errors.

For *object-action hallucination*, VideoHallu (Li et al., 2025c) uses synthetic “negative control” videos to test whether models incorrectly infer actions based on prior object–action associations instead of actual motion cues. For instance, a model may claim that a watermelon breaks after being shot even when it remains intact in the video. FactVC (Liu and Wan, 2023) identifies action consistency as a major source of captioning error, accounting for 38.3% of failures. Models often describe interactions such as a person dancing with a dog based on object co-occurrence, without grounding predictions in visual dynamics.

For *scene-event hallucination*, EventHallusion (Zhang et al., 2024) evaluates whether models hallucinate events by over-relying on typical scene–event pairings, such as assuming cooking takes place in a kitchen even without action evidence. NOAH (Lee et al., 2025) scales this evaluation to over 60,000 samples created from around 15,000 edited videos, testing whether models ignore inserted contradictory clips and instead generate events that align with the surrounding scene or narrative context. RoadSocial (Parikh et al., 2025)

focuses on driving scenarios, using adversarial and incompatible question formats to test whether models hallucinate common road events, such as collisions or traffic violations, based solely on general road context or misleading prompts, even when no such events occur in the video.

**Audio-visual conflict benchmarks** evaluate whether models integrate audio and visual signals appropriately, focusing on cases where dominant audio cues override visual input and lead to incorrect predictions. With only three existing benchmarks, this category remains underexplored, and current datasets are limited to short video clips. As multimodal Vid-LLMs increasingly process audio, further benchmark development is needed.

For *action attribution*, AVHBench (Sung-Bin et al., 2025) tests whether sounds such as music or bird calls cause models to generate incorrect visual descriptions like “a person is dancing” or “a bird is chirping,” even when no such actions are visible. It includes 2,136 videos and 5,302 binary question–answer pairs sourced from AudioCaps and VALOR, and reports precision, recall, and F1 scores to quantify errors. CMM (Leng et al., 2024) evaluates similar cases using curated “audio dominance” samples, where prominent sounds such as thunder occur without visual events like lightning. Models are asked binary questions to assess whether they mistakenly rely on audio alone for visual claims.

For *emotion inference*, EmotionHalluc (Xing et al., 2025) examines whether models infer incor-

rect emotional states based on misleading multi-modal cues. Its Reasoning Result and Reasoning Cue tasks test if models describe a neutral face as “excited” due to upbeat vocal tone, or invent emotional cues to justify unsupported conclusions.

### 4.3 Discussion: Coverage and Gaps

Table 1 summarizes 19 existing benchmarks for evaluating video hallucination, with a notable concentration on Spatiotemporal Dynamics (8 benchmarks), mostly targeting short clips. A few, such as Vript (Yang et al., 2024) and OVBench (Huang et al., 2025d), extend to long-form or streaming contexts. Context-Driven Fabrication shows broad domain coverage, ranging from traffic scenarios (RoadSocial) to synthetic videos (VideoHallu). In contrast, Referential Inconsistency and Audio-Visual Conflict remain underexplored, each represented by only three benchmarks, and no benchmark addresses audio-visual consistency in long-form videos. While 11 benchmarks include dedicated baselines to support method development, performance analysis reveals a clear divide: state-of-the-art models perform well on some tasks (e.g., VidHalluc, EventHallusion, with scores above 80%), but struggle with fine-grained temporal reasoning (e.g., VideoHalluc, 37.8%) and long-context consistency (e.g., ELV-Halluc, 53.1%). These findings identify dynamic distortion and long-range temporal grounding as persistent challenges for future research.

## 5 Mitigation Strategies

Following the taxonomy in Section 3, we group existing mitigation strategies by hallucination type and representative failure cases. Table 2 summarizes the corresponding techniques for each case.

### 5.1 Mitigating Dynamic Distortion

**Spatiotemporal dynamics mitigation** tackles event order, duration, and frequency Hallucinations using contrastive, optimization-based, and temporal grounding strategies, with event misordering being the most extensively studied.

For *event misordering*, SEASON (Wu et al., 2025a) contrasts original videos with temporally homogenized negatives that disrupt causal order, using self-diagnostic decoding to suppress outputs insensitive to correct sequence. Video-thinking (Gao et al., 2025) introduces TDPO (Thinking-based DPO), applying segment-weighted preference learning on reasoning paths

to optimize for temporal logic. SmartSight (Sun et al., 2025) ranks multiple responses based on the Temporal Attention Collapse (TAC) score, favoring outputs that attend proportionally across time to preserve correct order.

For *duration distortion*, Temporal Insight Enhancement (Sun et al., 2024) decomposes events into atomic actions and leverages external vision models to timestamp them, aligning model responses with grounded temporal claims. DINO-HEAL (Li et al., 2025a) uses DINOv2-guided spatial saliency to reweight features and maintain attention on action-relevant regions across time. TAAE (Cai et al., 2025) identifies activation offsets between full and downsampled inputs to amplify duration-sensitive representations during inference.

For *frequency confusion*, VTG-LLM (Guo et al., 2025) introduces absolute-time tokens that decouple event identity from repetition count, improving temporal anchoring of repeated actions. Vriptor (Yang et al., 2024) aligns dense scene-level captions with timestamps to enforce instance-level discrimination, helping models avoid merging or duplicating repeated events in long videos.

**Referential inconsistency mitigation** focuses on preserving distinct representations of entities and scenes across time.

For *character conflation*, Vista-LLaMA (Ma et al., 2024) introduces Equal Distance Attention, which removes positional decay between visual and textual tokens, ensuring stable attention to character identities regardless of when they appear. VideoPLR (Pu et al., 2025) constructs a structured video database with explicit object tracking, allowing symbolic logic programs to differentiate entity identities during reasoning.

For *scene conflation*, ELV-Halluc-DPO (Lu et al., 2025) applies adversarial preference optimization using cross-segment perturbations that swap entities in space or time, encouraging the model to ground predictions within the correct segment. VideoChat-Online (Huang et al., 2025d) uses a Pyramid Memory Bank during streaming inference to separate recent high-resolution frames from compressed long-term history, reducing confusion between distinct temporal contexts.

### 5.2 Mitigating Content Fabrication

**Context-driven fabrication mitigation** aims to separate predictions from statistical associations and strengthen visual grounding.

For *object-action hallucination*, methods em-

Table 2: Summary of video hallucination mitigation strategies. *Case*: **Order/Dur/Freq** (Spatiotemporal Dynamics), **Char/Scn** (Referential Inconsistency), **Obj-Act/Scn-Evt** (Context-Driven Fabrication), **Act/Emo** (Audio-Visual Conflict). *TF*:  $\checkmark$  = train-free,  $\times$  = training required. *Reported Gain*: reported improvement (over baseline) on primary benchmark.

Method	Venue	Case	TF	Core Technique	Key Mechanism	Reported Gain
<i>Spatiotemporal Dynamics Mitigation (Dynamic Distortion)</i>						
SEASON (Wu et al., 2025a)	arXiv'25	<b>Order</b>	$\checkmark$	Contrastive Decoding	Temporal homogenization contrast	+5.7% Acc (Qwen2.5-VL)
Video-thinking (Gao et al., 2025)	arXiv'25	<b>Order</b>	$\times$	Preference Optimization	Segment-weighted thinking contrast	+7.4% Acc (LLaVA-NeXT)
SmartSight (Sun et al., 2025)	arXiv'25	<b>Order</b>	$\checkmark$	Introspective Sampling	Temporal attention collapse score	+2.9% Acc (Video-R1)
Temporal Insight (Sun et al., 2024)	ICPR'24	<b>Dur</b>	$\checkmark$	Post-hoc Correction	Iconic action timestamp extraction	+27.9% R@1 (Video-LLaMA)
DINO-HEAL (Li et al., 2025a)	CVPR'25	<b>Dur</b>	$\checkmark$	Feature Reweighting	DINOv2 spatial saliency reweighting	+7.0% Acc (Video-LLaVA)
TAAE (Cai et al., 2025)	arXiv'25	<b>Dur</b>	$\times$	Activation Engineering	Temporal-aware offset injection	+4.4% Acc (Qwen2.5-VL)
VTG-LLM (Guo et al., 2025)	AAAI'25	<b>Freq</b>	$\times$	Temporal Grounding	Absolute-time token disentanglement	+6.8% R@1 (Video-LLaMA2)
Vriptor (Yang et al., 2024)	NeurIPS'24	<b>Freq</b>	$\times$	Video-Script Alignment	Dense script-based timestamp training	+7.5% F1 (ST-LLM)
<i>Referential Inconsistency Mitigation (Dynamic Distortion)</i>						
Vista-llama (Ma et al., 2024)	CVPR'24	<b>Char</b>	$\times$	Token Processing	Equal distance visual attention	~5.0% Acc (LLaVA)
VideoPLR (Pu et al., 2025)	arXiv'25	<b>Char</b>	$\times$	Perception-Logic-Reasoning	Database-anchored symbolic execution	+9.2% Acc (Qwen2.5-VL)
ELV-Halluc-DPO (Lu et al., 2025)	arXiv'25	<b>Scn</b>	$\times$	Preference Optimization	Cross-segment adversarial DPO	-27.7% SAH (Qwen2.5-VL)
VideoChat-Online (Huang et al., 2025d)	CVPR'25	<b>Scn</b>	$\times$	Streaming Processing	Pyramid memory bank update	+8.5% Acc (InternVL2)
<i>Context-Driven Fabrication Mitigation (Content Fabrication)</i>						
SANTA (Chang et al., 2025)	arXiv'25	<b>Obj-Act</b>	$\times$	Fine-grained Contrastive Tuning	Hard negative action/object swapping	+2.4% Acc (LLaVA-Video)
TCD (Zhang et al., 2024)	AAAI'26	<b>Obj-Act</b>	$\checkmark$	Contrastive Decoding	Logit subtraction of priors	+3.2% Acc (VILA)
MASH-VLM (Bae et al., 2025)	CVPR'25	<b>Scn-Evt</b>	$\times$	Disentangled Representation	DST-Attention & Harmonic-RoPE	+2.7% Acc (ST-LLM)
PaMi-VDPO (Ding et al., 2025)	arXiv'25	<b>Scn-Evt</b>	$\times$	Preference Optimization	Part-mismatch visual negatives	+5.9% Acc (LLaVA-OenVision)
MMA (Dang et al., 2025)	IJCAI'25	<b>Scn-Evt</b>	$\times$	Parameter-Efficient Tuning	Dual-path visual-text alignment	+2.0% Acc (MA-LMM)
VistaDPO (Huang et al., 2025a)	ICML'25	<b>O-A, S-E</b>	$\times$	Visual-State DPO	Penalizing low visual-dependency tokens	+36.5% Acc (Video-LLaVA)
VideoHallu-GRPO (Li et al., 2025c)	NeurIPS'25	<b>O-A, S-E</b>	$\times$	RL Fine-Tuning	Group-relative rewards on counter-intuitive data	+4.7% Acc (Qwen2.5-VL)
<i>Audio-Visual Conflict Mitigation (Content Fabrication)</i>						
AVHModel-Align-FT (Sung-Bin et al., 2025)	ICLR'25	<b>Act</b>	$\times$	Instruction Tuning	Modality-Disentangled Data	+33.8% Acc (Video-LLaMA)
AVCD (Jung et al., 2025)	NeurIPS'25	<b>Act</b>	$\checkmark$	Trimodal Contrastive Decoding	Dominance-aware Attentive Masking	+1.6% Acc (VideoLLaMA2)
PEP-MEK (Xing et al., 2025)	arXiv'25	<b>Emo</b>	$\checkmark$	Predict-Explain-Predict	Knowledge Extraction & Refinement	+9.1% Acc (Gemini2.5-Flash)

phasize motion cues over object presence. SANTA (Chang et al., 2025) applies fine-grained contrastive tuning using hard negatives where actions differ but entities remain the same, encouraging the model to rely on motion rather than co-occurrence patterns. TCD (Zhang et al., 2024) suppresses object-triggered predictions by subtracting logits from temporally shuffled inputs, guiding the model to attend to dynamic cues rather than static priors.

For *scene–event hallucination*, methods reduce the influence of background context on event inference. MASH-VLM (Bae et al., 2025) uses disentangled spatial-temporal attention to prevent the model from predicting actions based solely on static backgrounds. PaMi-VDPO (Ding et al., 2025) introduces preference learning with visually mismatched negatives to train the model to verify event descriptions against the actual scene. MMA (Dang et al., 2025) aligns local visual details with textual tokens through a dual-path adapter, reinforcing grounding in specific cues over general context.

Some methods address both cases. VistaDPO (Huang et al., 2025a) penalizes predictions that lack visual grounding by optimizing against prior-driven outputs, encouraging reliance on directly observable evidence. VideoHallu-GRPO (Li et al., 2025c) improves grounding using synthetic videos with counterintuitive scenarios, optimizing

model behavior through group-based relative rewards that favor visually consistent responses over learned priors.

**Audio-visual conflict mitigation** addresses errors caused by dominant audio signals overriding visual input. This area remains underexplored, with few targeted methods.

For *action attribution*, AVHModel-Align-FT (Sung-Bin et al., 2025) fine-tunes on annotations separating audio and visual events, helping models distinguish between auditory and visual sources. AVCD (Jung et al., 2025) uses contrastive learning with modality masking to suppress misleading cues and improve cross-modal grounding.

For *emotion inference*, PEP-MEK (Xing et al., 2025) enforces modality-specific reasoning by requiring models to explain visual evidence before integrating it with audio, reducing overreliance on vocal tone.

### 5.3 Discussion: Coverage and Trade-offs

Table 2 reveals uneven progress across hallucination types. While Spatiotemporal Dynamics and Context-Driven Fabrication are well addressed, Referential Inconsistency and Audio-Visual Conflict remain underexplored. A key trade-off exists between effectiveness and deployment cost. Training-based methods (e.g., VistaDPO (Huang et al., 2025a), AVHModel-Align-FT (Sung-Bin

524 *et al.*, 2025)) offer substantial gains (up to 30–36%)  
525 by reshaping model priors via reinforcement learn-  
526 ing or instruction tuning, but require high train-  
527 ing overhead. In contrast, training-free strategies  
528 (e.g., SEASON (Wu *et al.*, 2025a), SmartSight (Sun  
529 *et al.*, 2025), TCD (Zhang *et al.*, 2024)) are model-  
530 agnostic and easier to adopt, though with more  
531 limited gains. Mechanism suitability often aligns  
532 with error type: inference-time decoding or atten-  
533 tion adjustments help correct temporal and logical  
534 inconsistencies (e.g., VideoPLR (Pu *et al.*, 2025)),  
535 while hallucinations driven by strong priors call  
536 for deeper disentanglement during training (e.g.,  
537 MASH-VLM (Bae *et al.*, 2025)). Reducing the la-  
538 tency of current inference-time methods is critical  
539 for real-time and safety-sensitive deployment.

## 540 6 Future Directions

541 Building on the taxonomy in Section 3, we propose  
542 two core directions for enhancing hallucination ro-  
543 bustness in Vid-LLMs, targeting the underlying  
544 causes of dynamic distortion and content fabrica-  
545 tion.

### 546 6.1 Addressing Dynamic Distortion: 547 Temporal and Referential Fidelity

548 Dynamic distortion primarily results from a gap  
549 between visual encoding and temporal reasoning.  
550 This issue is often rooted in the use of static im-  
551 age encoders and pooling-based connectors, which  
552 tend to discard motion cues critical for capturing  
553 temporal dynamics (Li *et al.*, 2025b, 2023b; Zhang  
554 *et al.*, 2023b). As videos become longer, this prob-  
555 lem is further exacerbated by the limited capacity  
556 of current models to maintain long-range context,  
557 leading to semantic drift and referential inconsis-  
558 tency (Xiao *et al.*, 2024; Guo *et al.*, 2025).

559 To address these limitations, future architec-  
560 tures should adopt video-oriented designs that pre-  
561 serve temporal structure throughout the pipeline.  
562 This includes using video-native encoders such  
563 as VideoMAE (Wang *et al.*, 2023) and motion-  
564 aware connectors that incorporate signals like op-  
565 tical flow (Liu *et al.*, 2025b) to capture velocity  
566 and trajectory. For long-range consistency, models  
567 may benefit from structured memory mechanisms,  
568 including state space models like Mamba (Li *et al.*,  
569 2024) or episodic memory (Wang *et al.*, 2025), to  
570 retain persistent entity information over extended  
571 sequences.

### 572 6.2 Mitigating Content Fabrication: 573 Grounding and Alignment

574 Content fabrication arises when pretraining priors  
575 dominate over visual grounding. Models may hal-  
576 lucinate actions or events based on static entities or  
577 scenes, ignoring temporal evidence (Chang *et al.*,  
578 2025; Bae *et al.*, 2025). The problem is worsened  
579 by imbalanced modality integration, where dom-  
580 inant audio cues override visual input, leading to  
581 cross-modal conflicts (Sung-Bin *et al.*, 2025; Leng  
582 *et al.*, 2024).

583 To reduce fabrication, models should learn to  
584 separate priors from perceptual evidence. This  
585 can be achieved by using counterfactual strate-  
586 gies, such as introducing negative samples with  
587 implausible object–action pairs and applying debi-  
588 asing objectives that encourage reliance on motion  
589 cues (Qi *et al.*, 2024). In audio-visual settings,  
590 models should verify visual input before incorpo-  
591 rating audio signals to avoid hallucinations caused  
592 by sound (Mo and Song, 2024).

## 593 7 Conclusion

594 Video large language models (Vid-LLMs) have  
595 achieved significant progress in video-language  
596 modeling, but also give rise to unique hallucination  
597 patterns that differ from those in static image tasks.  
598 This survey introduces a mechanism-based taxon-  
599 omy that categorizes hallucinations in Vid-LLMs  
600 into two primary categories: Dynamic Distortion,  
601 referring to the misinterpretation of spatiotemporal  
602 progression or referential consistency; and Content  
603 Fabrication, referring to ungrounded outputs influ-  
604 enced by statistical context priors or dominant audi-  
605 tory cues. Although recent studies have advanced  
606 benchmarking and mitigation of spatiotemporal  
607 and context-driven hallucinations, challenges such  
608 as referential inconsistency and audio-visual con-  
609 flicts remain underexplored. Furthermore, most ex-  
610 isting mitigation strategies are applied at inference  
611 time or as post-training adjustments, highlighting  
612 the need for scalable, training-time alignment meth-  
613 ods. To improve the robustness of Vid-LLMs, we  
614 advocate future research toward developing video-  
615 native encoders that preserve motion cues, integrat-  
616 ing explicit memory mechanisms to support long-  
617 term temporal grounding, and employing counter-  
618 factual learning to disentangle model reasoning  
619 from prior-driven associations. These directions  
620 will be essential for building trustworthy and tem-  
621 porally faithful video-language systems.

## 622 Limitations

623 Previous surveys on hallucination in MLLMs have  
624 extended the scope from LLMs and image-based  
625 VLMs to include video and other modalities. How-  
626 ever, their coverage of video hallucination remains  
627 limited, with only brief mentions of benchmarks  
628 and mitigation efforts, lacking structured categor-  
629 ization or causal analysis. This survey addresses  
630 this gap by presenting a mechanism-driven tax-  
631 onomy of hallucinations in Vid-LLMs. We intro-  
632 duce a layered classification framework, review  
633 recent studies in greater depth, analyze the under-  
634 lying causes of hallucinations, and outline future  
635 research directions. While we have strived to cover  
636 key developments in Vid-LLM hallucination, some  
637 relevant work may be omitted. This survey includes  
638 research published up to January 2026.

## 639 Ethics Statement

640 This survey adheres to established ethical standards  
641 for academic research. All referenced works are  
642 publicly available, and no human subjects or per-  
643 sonally identifiable information are involved. The  
644 purpose of this survey is to facilitate academic un-  
645 derstanding and encourage responsible develop-  
646 ment in the study of hallucination in Vid-LLMs.  
647 All prior work has been properly cited, with appro-  
648 priate credit given to original contributions.

## 649 References

650 Kyungho Bae, Jinhyung Kim, Sihaeng Lee, Soonyoung  
651 Lee, Gunhee Lee, and Jinwoo Choi. 2025. MASH-  
652 VLM: mitigating action-scene hallucination in video-  
653 llms through disentangled spatial-temporal represen-  
654 tations. In *CVPR*, pages 13744–13753. Computer  
655 Vision Foundation / IEEE.

656 Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He,  
657 Zongbo Han, Zheng Zhang, and Mike Zheng Shou.  
658 2024. Hallucination of multimodal large language  
659 models: A survey. *CoRR*, abs/2404.18930.

660 Jianfeng Cai, Wengang Zhou, Zongmeng Zhang, Jiale  
661 Hong, Nianji Zhan, and Houqiang Li. 2025. Miti-  
662 gating hallucination in videollms via temporal-aware  
663 activation engineering. *CoRR*, abs/2505.12826.

664 Kai-Po Chang, Wei-Yuan Cheng, Chi-Pin Huang, Fu-En  
665 Yang, and Yu-Chiang Frank Wang. 2025. Mitigating  
666 object and action hallucinations in multimodal llms  
667 via self-augmented contrastive alignment. *CoRR*,  
668 abs/2512.04356.

669 Long Chen, Oleg Sinavski, Jan Hünemann, Alice Karn-  
670 sund, Andrew James Willmott, Danny Birch, Daniel

Maund, and Jamie Shotton. 2024. Driving with llms:  
Fusing object-level vector modality for explainable  
autonomous driving. In *ICRA*, pages 14093–14100. 671  
672 673

Wey Yeh Choong, Yangyang Guo, and Mohan S.  
Kankanhalli. 2024. Vidhal: Benchmarking temporal  
hallucinations in vision llms. *CoRR*, abs/2411.16771. 674  
675 676

Jisheng Dang, Shengjun Deng, Haochen Chang, Teng  
Wang, Bimei Wang, Shude Wang, Nannan Zhu, Guo  
Niu, Jingwen Zhao, and Jizhao Liu. 2025. Hallu-  
cination reduction in video-language models via hi-  
erarchical multimodal consistency. In *IJCAI*, pages  
9167–9175. ijcai.org. 677  
678 679 680 681 682

Xinpeng Ding, Kui Zhang, Jianhua Han, Lanqing Hong,  
Hang Xu, and Xiaomeng Li. 2025. Pami-vdpo:  
Mitigating video hallucinations by prompt-aware  
multi-instance video preference learning. *CoRR*,  
abs/2504.05810. 683  
684 685 686 687

Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen,  
Meng Zhao, Yifan Zhang, Xiong Wang, Di Yin, Long  
Ma, Xiawu Zheng, Ran He, Rongrong Ji, Yunsheng  
Wu, Caifeng Shan, and Xing Sun. 2024. VITA: to-  
wards open-source interactive omni multimodal LLM.  
*CoRR*, abs/2408.05211. 688  
689 690 691 692 693

Hongcheng Gao, Jiashu Qu, Jingyi Tang, Baolong  
Bi, Yue Liu, Hongyu Chen, Li Liang, Li Su, and  
Qingming Huang. 2025. Exploring hallucination  
of large multimodal models in video understand-  
ing: Benchmark, analysis and mitigation. *CoRR*,  
abs/2503.19622. 694  
695 696 697 698 699

Yongxin Guo, Jingyu Liu, Mingda Li, Dingxin Cheng,  
Xiaoying Tang, Dianbo Sui, Qingbin Liu, Xi Chen,  
and Kevin Zhao. 2025. VTG-LLM: integrating times-  
tamp knowledge into video llms for enhanced video  
temporal grounding. In *AAAI*, pages 3302–3310.  
AAAI Press. 700  
701 702 703 704 705

Haojian Huang, Haodong Chen, Shengqiong Wu, Meng  
Luo, Jinlan Fu, Xinya Du, Hanwang Zhang, and Hao  
Fei. 2025a. Vistadpo: Video hierarchical spatial-  
temporal direct preference optimization for large  
video models. In *ICML*. OpenReview.net. 706  
707 708 709 710

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong,  
Zhangyin Feng, Haotian Wang, Qianglong Chen,  
Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting  
Liu. 2025b. A survey on hallucination in large lan-  
guage models: Principles, taxonomy, challenges, and  
open questions. *ACM Trans. Inf. Syst.*, 43(2):42:1–  
42:55. 711  
712 713 714 715 716 717

Zhe Huang, Hao Wen, Aiming Hao, Bingze Song,  
Meiqi Wu, Jiahong Wu, Xiangxiang Chu, Sheng Lu,  
and Haoqian Wang. 2025c. Taming hallucinations:  
Boosting mllms’ video understanding via counterfac-  
tual video generation. *CoRR*, abs/2512.24271. 718  
719 720 721 722

Zhenpeng Huang, Xinhao Li, Jiaqi Li, Jing Wang, Xi-  
angyu Zeng, Cheng Liang, Tao Wu, Xi Chen, Liang  
Li, and Limin Wang. 2025d. Online video under-  
standing: Ovbench and videochat-online. In *CVPR*,  
723 724 725 726

727	pages 3328–3338. Computer Vision Foundation / IEEE.	780
728		781
729	Chaeyoung Jung, Youngjoon Jang, and Joon Son Chung.	782
730	2025. AVCD: mitigating hallucinations in audio-	783
731	visual large language models through contrastive de-	784
732	coding. <i>CoRR</i> , abs/2505.20862.	
733	Ming Kong, Xianzhou Zeng, Luyuan Chen, Yadong Li,	785
734	Bo Yan, and Qiang Zhu. 2025. Mhbench: Demysti-	786
735	fying motion hallucination in videollms. In <i>AAAI</i> ,	787
736	pages 4401–4409. AAAI Press.	788
737	Wei Lan, Wenyi Chen, Qingfeng Chen, Shirui Pan,	789
738	Huiyu Zhou, and Yi Pan. 2024. A survey of hal-	790
739	lucination in large visual language models. <i>CoRR</i> ,	791
740	abs/2410.15359.	792
741	Kyuhoo Lee, Euntae Kim, Jinwoo Choi, and Buru Chang.	793
742	2025. Noah: Benchmarking narrative prior driven	794
743	hallucination and omission in video large language	795
744	models. <i>arXiv preprint arXiv:2511.06475</i> .	796
745	Sicong Leng, Yun Xing, Zesen Cheng, Yang Zhou,	797
746	Hang Zhang, Xin Li, Deli Zhao, Shijian Lu, Chun-	798
747	yan Miao, and Lidong Bing. 2024. The curse of	799
748	multi-modalities: Evaluating hallucinations of large	800
749	multimodal models across language, visual, and au-	801
750	dio. <i>CoRR</i> , abs/2410.12787.	
751	Chaoyu Li, Eun Woo Im, and Pooyan Fazli. 2025a. Vid-	802
752	halluc: Evaluating temporal hallucinations in multi-	803
753	modal large language models for video understand-	804
754	ing. In <i>CVPR</i> , pages 13723–13733. Computer Vision	805
755	Foundation / IEEE.	806
756	Jinxuan Li, Chaolei Tan, Haoxuan Chen, Jianxin Ma,	807
757	Jian-Fang Hu, Wei-Shi Zheng, and Jianhuang Lai.	808
758	2025b. Image-to-video transfer learning based on	809
759	image-language foundation models: A comprehen-	810
760	sive survey. <i>CoRR</i> , abs/2510.10671.	811
761	Kunchang Li, Yanan He, Yi Wang, Yizhuo Li, Wen-	812
762	hai Wang, Ping Luo, Yali Wang, Limin Wang, and	813
763	Yu Qiao. 2023a. Videochat: Chat-centric video un-	814
764	derstanding. <i>CoRR</i> , abs/2305.06355.	815
765	Kunchang Li, Yanan He, Yi Wang, Yizhuo Li, Wen-	816
766	hai Wang, Ping Luo, Yali Wang, Limin Wang, and	817
767	Yu Qiao. 2023b. Videochat: Chat-centric video un-	818
768	derstanding. <i>CoRR</i> , abs/2305.06355.	819
769	Kunchang Li, Xinhao Li, Yi Wang, Yanan He, Yali	820
770	Wang, Limin Wang, and Yu Qiao. 2024. Video-	821
771	mamba: State space model for efficient video un-	822
772	derstanding. In <i>ECCV</i> , volume 15084, pages 237–255.	823
773	Springer.	824
774	Zongxia Li, Xiyang Wu, Guangyao Shi, Yubin Qin,	825
775	Hongyang Du, Tianyi Zhou, Dinesh Manocha,	826
776	and Jordan Lee Boyd-Graber. 2025c. Videohallu:	827
777	Evaluating and mitigating multi-modal hallucina-	828
778	tions on synthetic video understanding. <i>CoRR</i> ,	829
779	abs/2505.01481.	830
	Bin Lin, Yang Ye, Bin Zhu, Jiayi Cui, Munan Ning,	831
	Peng Jin, and Li Yuan. 2024. Video-llava: Learning	832
	united visual representation by alignment before pro-	833
	jection. In <i>EMNLP</i> , pages 5971–5984. Association	834
	for Computational Linguistics.	835
	Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen,	
	Xiutian Zhao, Ke Wang, Liping Hou, Rongjun Li, and	
	Wei Peng. 2024a. A survey on hallucination in large	
	vision-language models. <i>CoRR</i> , abs/2402.00253.	
	Hui Liu and Xiaojun Wan. 2023. Models see halluci-	
	nations: Evaluating the factuality in video caption-	
	ing. In <i>EMNLP</i> , pages 11807–11823. Association	
	for Computational Linguistics.	
	Jing Liu, Sihan Chen, Xingjian He, Longteng Guo,	
	Xinxin Zhu, Weining Wang, and Jinhui Tang. 2025a.	
	VALOR: vision-audio-language omni-perception pre-	
	training model and dataset. <i>IEEE Trans. Pattern Anal.</i>	
	<i>Mach. Intell.</i> , 47(2):708–724.	
	Ruyang Liu, Shangkun Sun, Haoran Tang, Wei Gao,	
	and Ge Li. 2025b. Flow4agent: Long-form video	
	understanding via motion prior from optical flow. In	
	<i>CVPR</i> , pages 23817–23827.	
	Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang,	
	Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and	
	Lu Hou. 2024b. Tempcompass: Do video llms really	
	understand videos? In <i>ACL (Findings)</i> , pages 8731–	
	8772. Association for Computational Linguistics.	
	Hao Lu, Jiahao Wang, Yaolun Zhang, Ruohui Wang,	
	Xuanyu Zheng, Yepeng Tang, Dahua Lin, and Lewei	
	Lu. 2025. Elv-halluc: Benchmarking semantic ag-	
	gregation hallucinations in long video understanding.	
	<i>CoRR</i> , abs/2508.21496.	
	Fan Ma, Xiaojie Jin, Heng Wang, Yuchen Xian, Jiashi	
	Feng, and Yi Yang. 2024. Vista-llama: Reducing	
	hallucination in video language models via equal	
	distance to visual tokens. In <i>CVPR</i> , pages 13151–	
	13160. IEEE.	
	Muhammad Maaz, Hanoona Abdul Rasheed, Salman	
	Khan, and Fahad Khan. 2024. Video-chatgpt: To-	
	wards detailed video understanding via large vision	
	and language models. In <i>ACL (I)</i> , pages 12585–	
	12602. Association for Computational Linguistics.	
	Shentong Mo and Yibing Song. 2024. Aligning audio-	
	visual joint representations with an agentic workflow.	
	In <i>NeurIPS</i> .	
	Chirag Parikh, Deepti Rawat, Rakshitha R. T, Tatha-	
	gata Ghosh, and Ravi Kiran Sarvadevabhatla. 2025.	
	Roadsocial: A diverse videoqa dataset and bench-	
	mark for road event understanding from social video	
	narratives. In <i>CVPR</i> , pages 19002–19011. Computer	
	Vision Foundation / IEEE.	
	Bowei Pu, Chuanbin Liu, Yifan Ge, Peichen Zhou,	
	Yiwei Sun, Zhiyin Lu, Jiankang Wang, and Hong-	
	tao Xie. 2025. Alternating perception-reasoning for	
	hallucination-resistant video understanding. <i>CoRR</i> ,	
	abs/2511.18463.	



- 944 Jiacheng Zhang, Yang Jiao, Shaoxiang Chen, Jingjing  
945 Chen, and Yu-Gang Jiang. 2024. Eventhallusion:  
946 Diagnosing event hallucinations in video llms. *CoRR*,  
947 abs/2409.16597.
- 948 Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu,  
949 Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang,  
950 Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei  
951 Bi, Freda Shi, and Shuming Shi. 2023c. Siren’s song  
952 in the AI ocean: A survey on hallucination in large  
953 language models. *CoRR*, abs/2309.01219.
- 954 Fufangchen Zhao, Liao Zhang, Daiqi Shi, Yuanjun Gao,  
955 Chen Ye, Yang Cai, Jian Gao, and Danfeng Yan.  
956 2025. Videoperceiver: Enhancing fine-grained tem-  
957 poral perception in video multimodal large language  
958 models. *CoRR*, abs/2511.18823.
- 959 Luowei Zhou, Chenliang Xu, and Jason J. Corso. 2018.  
960 Towards automatic learning of procedures from web  
961 instructional videos. In *AAAI*, pages 7590–7598.  
962 AAAI Press.