

Figure 2: The overall framework of MLLM Alignment

tailed technical pipeline for multimodal alignment, covering preference signal construction and alignment objective design, thereby formalizing an end-to-end recipe for the MLLM alignment algorithm.

- **Application scenarios.** We analyze how alignment manifests across a variety of multimodal application scenarios, highlighting recurring design patterns as well as task-dependent requirements and constraints.
- **Evaluation framework.** Across diverse application scenarios, we systematically compare the relevant evaluation benchmarks and their key emphases, thereby guiding researchers to refine alignment strategies toward metric-specific objectives correspondingly.

The overall framework of our paper is shown

in Fig 2. This survey aims to provide a comprehensive and systematized perspective that enables researchers in both academia and industry to identify fundamental, cross-domain challenges in multimodal alignment. By doing so, we hope it can facilitate the design and innovation of more broadly applicable alignment algorithms that can effectively address, and ultimately go beyond, the complexities of real-world environments.

2 MLLM Alignment Technical Pipeline

Alignment in MLLM is generally an end-to-end design process. In practice, it is typically decomposed into two major stages: **preference signal construction** and **alignment objective design**. The preference signal determines the sources and representational forms of supervision used during alignment. In contrast, the alignment objective specifies how these signals are integrated into learning and the

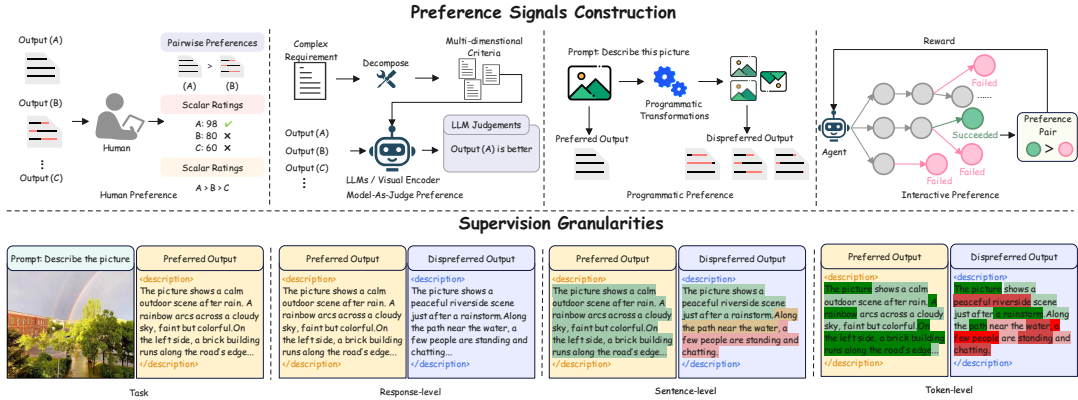


Figure 3: Preference signal construction mechanisms and supervision granularity designs.

degree or intensity of alignment is enforced. Consequently, alignment across domains and tasks should not be viewed as isolated cases; rather, it can be understood as an adaptive design grounded in diverse human preferences and expectations. This section provides a detailed decomposition and analysis of this unified design framework.

2.1 Preference Signals Construction

Multimodal settings contain rich, experience-grounded informational factors. Expected preferences vary substantially across tasks and may even be mutually conflicting. This diversity gives rise to heterogeneous, task-specific alignment needs, which in turn impose different requirements on how preference signals should be constructed, as shown in the upper part of Fig 3.

2.1.1 Human Preference

The most direct form of supervision for preference alignment comes from humans. To a considerable extent, MLLM alignment mirrored the success of RLHF in the text-only setting: annotators provide either pairwise preferences, scalar ratings, or ranked sequences to support preference labeling. Early MLLMs alignment methods, such as LLaVA-RLHF (Sun et al., 2023a), RLHF-V (Yu et al., 2024a), M-HalDetect (Gunjal et al., 2024), largely relied on carefully designed human annotations for training datasets, which substantially mitigated hallucinations and jailbreak attacks. In domains with high professional complexity or stringent safety requirements, human-in-the-loop remains indispensable. In biomedical systems, 3D-CT-GPT++ (Chen et al., 2025a) depend on expert-designed tasks and quality criteria, even when parts of the subsequent

scoring process are delegated to language models. Similarly, embodied settings such as EMMOE (Li et al., 2025a) typically require the collection and analysis of expert demonstrations as a prerequisite for building imitation-learning and preference datasets.

2.1.2 Model-As-Judge Preference

Golden standard human annotations are prohibitively expensive. With the growing adoption of RLAIIF (Lee et al., 2024b), increasing work has shifted to the LLM-as-judge paradigm, leveraging strong LLMs or auxiliary visual encoders to generate labels and scores for candidate samples. Closed-source models such as GPT (OpenAI, 2024; Achiam et al., 2023) and Gemini (Comanici et al., 2025; Team, 2024) have reached near-human performance on multiple tasks, therefore widely used to annotate preference data across diverse areas, including vision-language understanding (Liu et al., 2023b; Li et al., 2023a; Lee et al., 2024c), video understanding (Tang et al., 2025; Sun et al., 2024), long chain-of-thought reasoning (Zhang et al., 2025b; Wang et al., 2025e), and model safety (Zong et al., 2024; Zhao et al., 2024).

A key ingredient of Model-based annotation is prompt engineering that operationalizes target preferences: decomposing complex requirements into multi-dimensional criteria to improve reliability. For example, Silkie (Li et al., 2023a) uses GPT-4V to score responses along dimensions such as helpfulness, visual faithfulness, and ethical considerations, yielding relatively stable preference signals. CLIP-DPO (Ouali et al., 2024), in contrast, employs CLIP-like visual encoders to assist preference generation and scoring, illustrating how

171	modular encoders can be exploited for preference	219
172	construction.	
173	Beyond such external judges, preference sig-	220
174	nals could come from the target model itself: de-	221
175	pending on modality settings, methods such as	222
176	SQuBa (Eom et al., 2025), SymDPO (Jia et al.,	223
177	2024), SIMA (Wang et al., 2024e), and MIA-	224
178	DPO (Liu et al., 2024c) construct preference	225
179	pairs from model responses in text-only or vision-	226
180	language scenarios. Nevertheless, Model-as-judge	227
181	substantially improves labeling efficiency and re-	228
182	duces cost, but it also inevitably transfers the biases	229
183	and failure modes of the judging or generating mod-	230
184	els into the preference data.	231
185	2.1.3 Programmatic Preference	232
186	Commonly, preferences are programmatically spec-	233
187	ified through reusable rules and procedural gener-	234
188	ation. This makes programmatic preference con-	235
189	struction especially important for MLLMs, as it	236
190	directly encodes task structure, safety constraints,	237
191	and definitions of hallucination into the data gener-	238
192	ation pipeline.	239
193	Visual perturbations and contrastive variants pro-	240
194	vide a practical route to constructing preference	241
195	data at scale. Image DPO (Luo et al., 2025) and	242
196	AdPO (Liu et al., 2024a) generate aligned posi-	243
197	tive/negative samples through semantic edits, blur-	244
198	ring, or adversarial transformations, while MIA-	245
199	DPO and PanoDPO (Liang et al., 2025a) exploit	246
200	compositing, viewpoint variation, or conditional	247
201	preferences to strengthen visual grounding. For	248
202	hallucination and safety control, HDPO (Fu et al.,	249
203	2024c), CcDPO (Li et al., 2025b), SPR (Qiu et al.,	250
204	2025), and SAFEVID (Wang et al., 2025f) create	251
205	contrastive pairs by modulating hallucination sever-	252
206	ity, temporal coherence, or exposure to hazardous	253
207	content. Moving beyond strict pairwise setups,	254
208	MCM-DPO (Fu et al., 2025c) and MFPO (Jiang	255
209	et al., 2025a) define modality-aware, multi-criteria	256
210	preferences, and MMedPO (Zhu et al., 2025a) in-	257
211	troduces domain-aware weighting to emphasize	258
212	clinically consequential errors.	259
213	Programmatic preference design has become	260
214	central to modern MLLM alignment and offers	261
215	low-cost scalability, improved control over specific	262
216	failure modes, and more reproducible training sig-	263
217	nals than large-scale human annotation or generic	264
218	judge-based scoring.	265
	2.1.4 Interactive Preference	266
	As the scope of MLLM tasks expands, preference	267
	signals no longer arise solely from static annotation	268
	pipelines; they can also be induced from reward	
	information obtained through interaction with envi-	
	ronments or tasks. MLLMs operate as agents that	
	generate actions, queries, or reasoning trajectories	
	and receive feedback in the form of task success sig-	
	nals, environment rewards, or trajectory-level per-	
	formance scores. INTERACTIVECOT (Jiao et al.,	
	2025) in ALFWorld, EMMOE in Replica-like 3D	
	environments, and D ² PO (Wang et al., 2025b) for	
	embodied planning, which derive preferences from	
	task completion signals and decomposed subgoals.	
	GRAPE similarly aligns policies at the trajectory	
	level by contrasting successful and failed trials un-	
	der spatiotemporal constraints.	
	While auxiliary LLM judges may still be em-	
	ployed at certain stages, the primary supervision	
	in these approaches is inherently dynamic and	
	trajectory-dependent rather than a one-shot evalua-	
	tion of static input–output pairs, making interaction	
	itself a central source of alignment feedback.	
	2.2 Alignment Objective Design	
	Alignment objectives for multimodal build upon	
	the general framework of preference-based or	
	feedback-based alignment, while introducing	
	broader constraints and more fine-grained de-	
	compositions centered on optimization paradigms,	
	structural granularity and stability mechanisms.	
	The detailed optimization objectives are shown in	
	Table 1.	
	2.2.1 Optimization Paradigms	
	Online optimization largely follows the RLHF-	
	style formulation, treating alignment as a policy-	
	based reinforcement learning problem. This was a	
	commonly adopted framework in early alignment	
	work; Fact-RLHF (Sun et al., 2023a) uses human	
	scores as rewards to construct the RLHF training	
	process. Related RL variants, such as DDPO (Yu	
	et al., 2024a) and RL4VLM (Zhai et al., 2024) in-	
	teract with an environment or sampler and receive	
	reward signals provided by humans or strong base-	
	line models.	
	Instead of explicitly fitting a reward model and	
	performing on-policy policy optimization as in	
	PPO-based RLHF, offline-based DPO directly opti-	
	mizes the policy with a simple contrastive objective	
	over offline preference pairs, while implicitly in-	
	corporating a reference-model regularization effect.	

This RL-free formulation substantially reduces algorithmic and engineering complexity, improves training stability, and lowers computational cost by avoiding iterative sampling-and-optimization loops. Together with its simplicity and scalability, it has made it one of the most widely adopted approaches for preference alignment in current practice.

2.2.2 Structural Granularity

Traditional alignment typically operates at the response level, treating the output as a single response sequence. For a multimodal input (x, \mathcal{I}) , where x is a textual query and \mathcal{I} denotes the associated images. The objective is commonly defined over preference pairs, optimizing the model to favor a preferred response over a dispreferred response (i.e., (y_w, y_l)). Although this response-level formulation is effective across many tasks, recent works move beyond it by introducing fine-grained and hierarchical supervision illustrated in Fig 3 below, either by assigning non-uniform weights to tokens/spans or by defining dedicated objectives for specific structured components of multimodal outputs.

Fine-grained sequence. Fine-grained supervision improves credit assignment by pushing preference signals down to segment, sentence, or token units, which helps suppress localized hallucinations and makes better use of limited preference data. Dense-DPO in RLHF-V concentrates optimization on uncertain or reliability-critical spans, so the model is explicitly trained to correct the most error-prone fragments rather than merely shifting overall response style. FDPO (Gunjal et al., 2024) refines the pairwise objective by decoupling the contributions inside the sigmoid term, offering tighter control over local gradient behavior while keeping the offline preference optimization simplicity. This line also motivates token-level or span-aware DPO variants such as fdPO (Shen et al., 2025) and VAD-DPO (Zhang et al., 2025a) that retain the response-level backbone but inject finer error-targeting mechanisms.

Multiple hierarchies. Hierarchical fusion improves stability and controllability by jointly aligning global preference direction and local grounding constraints across different granularities. HA-DPO (Zhao et al., 2024) and MIA-DPO extend response-level DPO with auxiliary supervised losses from positive examples to strengthen alignment under noisy or sparse supervision. CHiP (Fu

et al., 2025c) represents a more explicit multi-level design by combining visual DPO, sentence-level DPO, and a token-level KL-style sequence constraint, effectively distributing alignment pressure across textual structure and visual evidence. Even when the main objective remains response-level, MIA-DPO’s attention-based multi-image pairing implicitly spreads supervision over multiple visual components, improving robustness to compositional inputs.

Multiple dimension. Multi-dimensional designs broaden alignment beyond a single “good vs bad” axis, enabling structured control over multi-component inputs and heterogeneous error modes. MCM-DPO and MFPO define preferences over combinations of answer text, images, and context, encouraging consistent cross-modal reasoning rather than modality-isolated optimization. MMedPO and HSCR (Jiang et al., 2025b) introduce multiple dispreference types and learn to rank them, which supports finer-grained governance of distinct failure patterns such as hallucination versus omission in domain-critical settings. Video-SALMONN series (Sun et al., 2024; Tang et al., 2025) extends this idea to temporal media, video-oriented objectives treat coherence and time-consistency as additional alignment dimensions, pushing preference learning toward richer, scenario-aware control.

2.2.3 Regularization and Stability

As many alignment algorithms can be formulated as contrastive preference optimization, a key concern is preventing excessive distribution shift, catastrophic forgetting, and undesirable likelihood drift relative to base models. Recent work addresses this with two recurring objective-level designs.

Reference-oriented regularization. Reference-oriented regularization constrains updates toward a trusted baseline. Fact-RLHF uses an explicit KL penalty to limit PPO deviation from the initial policy; CHiP adds KL terms between the reference and updated models for both preferred and dispreferred responses; and PO/MPO-style objectives incorporate log-ratio or reference-based penalties to keep adapted policies close to clean or accurate anchors.

Margin and asymmetry shaping. Complementing reference constraints, margin design and asymmetry-aware objectives regulate how strongly preferences are enforced. DAMA (Lu et al., 2025) rescales the DPO logit gap, mDPO (Wang et al.,

2024a) and MPO (Wang et al., 2024d) impose margin-style constraints to preserve base accuracy while enforcing a minimum preference distance, and hallucination-oriented variants such as HDPO and SymDPO adopt anchored or symmetric formulations to mitigate probability drift.

3 Application Scenarios

MLLM alignment is inherently application-driven. As the scope of MLLM deployment continues to expand, these models have demonstrated substantial value in both general-purpose multimodal settings and more domain-specific scenarios.

3.1 General Multimodality

General multimodal capability serves as a foundational application for assessing the effectiveness of multimodal alignment. Correspondingly, at the level of alignment strategy, it manifests as four system-level dimensions that recur across multimodal model designs: hallucination mitigation, instruction following, reasoning consistency, and robustness and resistance.

Hallucination mitigation Hallucination mitigation is a core objective in MLLM alignment, aiming to ensure responses are grounded in visual evidence rather than linguistic priors or spurious cues. Most methods construct preference pairs that differ in visual faithfulness and apply token- or instance-level penalties. Clause-aware approaches (Zhao et al., 2024) suppress labeled hallucinated spans; image-perturbation methods such as CLIP-DPO, Image DPO and AdPO enforce true visual dependence by altering visuals while keeping text fixed; and multi-image PanoDPO (Liang et al., 2025a) and LPOI (Zadeh et al., 2025) contrast viewpoints or occlusions to reduce shortcut reasoning. The inclusion of visual faithfulness as an explicit scoring in broader systems (Lee et al., 2024a) indicates that hallucination control is now a standard component of general MLLM alignment.

Instruction following Instruction-following alignment aims to position MLLMs as practical, well-behaved assistants that can correctly interpret complex prompts, respect constraints, and produce task-compliant outputs. RLHF-style pipelines optimize policies with user-derived rewards that reflect overall user preferences over multimodal responses. In contrast, DPO-style (Liu et al., 2025a; Xue et al., 2024; Li et al., 2024b; Xiong

et al., 2024) scales instruction-response data construction and derives preference pairs from human ratings, LLM-as-judge scoring, or cross-model comparisons. Extensions such as FDPO, DAMA, and mDPO refine how strong versus weak preferences are separated, while RLAIIF-V (Yu et al., 2024b) and RLHF-V incorporate dense judgments from proprietary or open evaluators to improve controllability and reliability. Across these systems, visual inputs are integrated into prompts, yet the primary alignment target remains robust multimodal instruction adherence and consistent, controllable satisfaction of user goals.

Reasoning consistency With the rise of CoT, multimodal reasoning has become an increasingly active direction in MLLM alignment. Reasoning alignment targets how models reach answers by shaping CoT traces, decomposition, and intermediate decisions rather than only final outputs. This line spans domain-focused CoT alignment for math and visual reasoning (Zhang et al., 2024a,b), step-wise optimization for temporal and audio-visual understanding (Tang et al., 2024; Sun et al., 2025), and environment-grounded CoT improved by task success signals in embodied settings (Jiao et al., 2025). In parallel, several methods align multi-step intermediate artifacts, such as plans (Song et al., 2025), code variants (Zhang et al., 2025c), or structured edits (Zhu et al., 2025b; He et al., 2025), so that preference learning directly constrains the reasoning products themselves. These approaches broaden reasoning alignment from explicit CoT supervision to outcome-centric and artifact-centric designs, encouraging MLLMs to act as more interpretable, human-aligned multimodal reasoners.

Robustness and resistance Robustness and resistance alignment aims to ensure that MLLMs remain stable under benign perturbations and resilient to adversarial or misleading inputs, such as inconsistencies across images or between modalities, thereby reducing modality inertia and spurious cross-modal reliance. Image DPO and AdPO construct preference pairs by perturbing only the visual input, rewarding responses that remain correct when semantics are preserved and appropriately change when visual meaning shifts. VAD-DPO extends this idea to videos by contrasting visually similar yet semantically opposite clips to heighten sensitivity to subtle temporal differences. Complementary robustness-oriented objectives, such as SymDPO, SymMPO (Liu et al., 2025b), adopt sym-

467	metric, to curb over-reliance on textual cues and	516
468	preserve performance. While MCM-DPO, MFPO	517
469	and HSCR use sparsity-aware formulations that	518
470	generalize supervision from isolated pairs to multi-	519
471	candidate comparisons and encourage reliance on	520
472	the most informative visual evidence.	521
473	3.2 Domain-specific Tasks	522
474	MLLM alignment becomes more in-depth in	523
475	domain-specific tasks, exhibiting strong domain	524
476	specificity, and is primarily applied to Medical and	525
477	science, embodied AI, and safety scenarios.	526
478	Medical and science Medical and scientific ap-	527
479	plications place a premium on factual correctness,	528
480	domain expertise, and calibrated behavior under	529
481	complex multimodal evidence. 3D-CT-GPT++ con-	530
482	structs radiology preference data scored by ex-	531
483	perts or experienced reasoning LLMs. MMedPO	532
484	differentiates hallucination-type versus lesion and	533
485	weights them by clinical severity with visualization-	534
486	aware signals. MSR-ViR further aligns long-form	535
487	scientific summarization via preferences over multi-	536
488	stage outputs. HSCR (Jiang et al., 2025b) targets	537
489	medical VLM alignment by automatically con-	538
490	structing high-quality preference pairs to improve	539
491	clinical trustworthiness with limited training data.	540
492	Embodied AI MLLMs are increasingly used for	541
493	spatial perception and action selection in simulated	542
494	and real environments. INTERACTIVECOT col-	543
495	lects rollouts in environments like ALFWorld and	544
496	treats task success as a preference signal to opti-	545
497	mize multimodal CoT traces for action. In robotics,	546
498	GRAPE aligns manipulation policies by contrast-	547
499	ing successful and failed trajectories under spatial	548
500	constraints, while D ² PO jointly learns action selec-	549
501	tion and state prediction via preference learning,	550
502	coupling policy and world model learning. Domain	551
503	systems like AIGI-Holmes (Zhou et al., 2025) com-	552
504	bine vision experts with MLLMs and apply DPO	553
505	to task-specific datasets for visual forensics and in-	554
506	spection. Across these settings, alignment depends	555
507	on environmental dynamics and long-horizon suc-	556
508	cess, making robustness to distribution shift and	557
509	stable interaction crucial.	558
510	Safety Safety-centered alignment seeks to pre-	559
511	vent harmful content, curb misuse, and improve	560
512	risk awareness in multimodal systems. BPO (Li	561
513	et al., 2025c) introduces the MMSafe-PO dataset	562
514	and a text-inertia loss to deter “modality cheat-	563
515	ing,” making safe, visually grounded responses	564
	strictly preferred over unsafe or ungrounded ones.	
	Training-time safety alignment is increasingly com-	
	plemented by inference-time modules that serve	
	as multimodal filters or risk classifiers. LLaMA	
	Guard 3 Vision (Cheng et al., 2024) provides pair-	
	wise safety gating at the response level, while VL-	
	Guard (Zong et al., 2024) targets lightweight guard	
	models and calibrated risk detection for VLM de-	
	ployments. MultiTrust (Zhang et al.) offers a	
	unified trustworthiness evaluation suite to audit	
	safety and robustness under diverse stressors, and	
	MMDT (Xu et al., 2025) enforces decoding-time	
	trustworthiness by intervening during generation	
	rather than relying only on post-hoc filtering. Over-	
	all, safety alignment typically adopts conservative	
	objectives and uses adversarially constructed pref-	
	erence data to improve robustness and reliability	
	against hostile or deceptive inputs.	
	4 Benchmark and Evaluations	
	Mirroring the application landscape, existing	
	benchmarks can be broadly organized into general	
	multimodal and domain-specific evaluations: the	
	former are designed to reflect breadth-oriented use	
	of MLLMs across diverse vision–language tasks	
	and thus measure average capability and overall	
	alignment quality, while the latter track deploy-	
	ment in specialized settings with task-critical depth	
	and boundary conditions. We list representative	
	benchmarks in Table 2.	
	4.1 General Multimodal Benchmarks	
	General-purpose multimodal benchmarks are typi-	
	cally designed to probe broad, cross-task capabili-	
	ties over diverse vision–language settings, a trend	
	reinforced by the growing integration of alignment	
	into mainstream post-training pipelines. The focus	
	of general evaluation has progressively crystallized	
	into three dimensions: knowledge and reasoning,	
	hallucination and instruction, and unified prefer-	
	ence assessment.	
	Knowledge and reasoning. As demands for mul-	
	timodal reasoning grow, general-purpose knowl-	
	edge and reasoning benchmarks are increasingly	
	used to track how alignment reshapes a model’s	
	overall reasoning profile. MMMU (Yue et al.,	
	2024) and MathVista (Lu et al., 2023) remain	
	representative suites for cross-disciplinary and	
	math-oriented reasoning, while SQA3D (Ma et al.,	
	2023) targets spatial perception and grounded rea-	
	soning. Complementing these, fine-grained real-	

565	world visual reasoning benchmarks are often orga-	expert-defined criteria rather than broad average	615
566	nized as broader suites (Jiang et al., 2024; Zhang	performance.	616
567	et al., 2024d; Chen et al., 2024a; Liu et al., 2023c;		
568	Ying et al., 2024; Fu et al., 2024b), which probe	Medical and science. Evaluations in the medical	617
569	more localized and realistic visual reasoning effects	and science fields prioritize expert-level accuracy	618
570	of alignment.	and norms, characterized by an extremely low tol-	619
571	Hallucination and instruction. Another widely	erance for error. Med-VQA (Canepa et al., 2023)	620
572	emphasized dimension is whether aligned mod-	and 3D-CT-GPT++ quantify proficiency in inter-	621
573	els can produce faithful, evidence-consistent re-	preting specialized imaging, whereas MMedPO	622
574	sponses and remain robust against misleading or	and related methodologies define criteria for clini-	623
575	shifting contexts. Benchmarks such as MMHal-	cal preference alignment. Chen et al. ensure that	624
576	Bench (Sun et al., 2023b), VHTest (Huang et al.,	model outputs strictly adhere to diagnostic logic	625
577	2024), GAVIE (Liu et al., 2023a), and HQH (Yan	and terminological standards. In immune adjuvants	626
578	et al., 2024) primarily measure factual ground-	and materials chemistry, expert-designed standards	627
579	ing by testing consistency between visual inputs and	and consistency validation were introduced to en-	628
580	generated responses, enabling systematic measure-	sure the reliability of assessments.	629
581	ment of multimodal hallucinations. In contrast, R-		
582	Bench (Wu et al., 2024), VLind-Bench (Lee et al.,	Embodied AI. Benchmarks in the Embodied AI	630
583	2024c), Vibe-Eval (Padlewski et al., 2024), and	domain aim to validate closed-loop "perception-	631
584	LiveBench (White et al., 2024) place greater weight	decision-action" capabilities, facilitating the tran-	632
585	on robustness and interaction behavior, assessing	sition from static visual processing to dynamic in-	633
586	whether models maintain cautious, stable, and faith-	teraction. While SQA3D establishes rigorous stan-	634
587	ful response patterns under adverse conditions or	dards for 3D spatial perception involving depth	635
588	distribution shifts.	and occlusion, ASTRA (Wang et al., 2025a) and	636
589	Unified preference assessment. Alignment	Engine scrutinize task planning and navigation by	637
590	stems from the need for consistency with human	translating high-level instructions into executable	638
591	preferences, so a model’s alignment is often judged	actions, and the RL4VLM framework further val-	639
592	by how closely its outputs match human judgments.	idates reinforcement learning strategies based on	640
593	Q-Bench (Wu et al., 2023), LLVisionQA, LLDe-	environmental feedback.	641
594	scribe, and LLaVA Bench-Wilder (Li et al., 2024a)		
595	evaluate alignment to human expectations via	Safety. Safety-oriented benchmarks prioritize	642
596	controlled perception-focused QA and description	stress-testing model defenses to construct robust	643
597	tasks, jointly reflecting visual ability and instruc-	adversarial protections and define the boundaries	644
598	tion following. Going beyond task performance,	of harmlessness. VLGuard and VLSBench (Hu	645
599	M-RewardBench (Gureja et al., 2024), MJ-	et al., 2024) quantify the ability to reject "Visual	646
600	Bench (Chen et al., 2024b), VL-RewardBench (Li	Jailbreak" attacks; RLAIIF-V and HallusionBench	647
601	et al., 2024c), and RewardBench (Lambert et al.,	focus on mitigating harmful multimodal halluci-	648
602	2024) assess alignment at the preference level,	nations; and SafeVID extends these ethical compli-	649
603	measuring how well multimodal reward models	ance assessments to dynamic video generation.	650
604	or LLM-as-judge systems approximate human		
605	preferences on safety, reliability, and bias in	5 Conclusion	651
606	multi-task, multilingual settings.		
607	4.2 Domain-specific Benchmarks	MLLM alignment has progressed rapidly in recent	652
608	Domain-specific benchmarks narrow the evaluation	years. In this survey, we provide a systematic and	653
609	scope and pose a more explicit question: Does the	comprehensive review of existing work from three	654
610	model meet the acceptable and useful standards	complementary perspectives: technical pipelines,	655
611	of a given domain? These evaluations are inher-	application scenarios, and evaluation methodolo-	656
612	ently more domain-sensitive, typically grounded	gies. Due to space constraints, we place an ex-	657
613	in concrete tasks and downstream use cases, and	tended discussion of future directions in the Ap-	658
614	emphasize boundary conditions, risk profiles, and	pendix D. We hope this survey will serve as a useful	659
		reference and help catalyze further research on prin-	660
		cipled, reliable, and scalable MLLM alignment.	661

662 Limitations

663 The paper retrieval, inclusion, and exclusion processes were performed by a single reviewer (the
664 study’s first author). While we implemented rigorous procedures to ensure comprehensive coverage
665 of published works, this approach inherently carries the risk of omitting potentially relevant studies.
666 Furthermore, classification of papers into specific categories or citation implementation might contain
667 inadvertent errors. Nevertheless, we have performed multiple verification steps throughout the
668 analytical process to mitigate such limitations. Although minor inconsistencies or omissions may
669 persist, we maintain that this survey constitutes the most comprehensive review of MLLM alignment
670 currently available, offering an objective and detailed assessment of future research directions and
671 outstanding challenges.
672
673
674
675
676
677
678
679

680 References

681 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
682 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
683 Diogo Almeida, Janko Altenschmidt, Sam Altman,
684 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.
685 *arXiv:2303.08774*.

686 Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna,
687 Baptiste Bout, Devendra Chaplot, Jessica Chud-
688 novsky, Diogo Costa, Baudouin De Monicault,
689 Saurabh Garg, Theophile Gervet, et al. 2024. Pixtral
690 12b. *arXiv preprint arXiv:2410.07073*.

691 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang,
692 Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou,
693 and Jingren Zhou. 2023. Qwen-vl: A frontier large
694 vision-language model with versatile abilities. *arXiv
695 preprint arXiv:2308.12966*.

696 Shuai Bai et al. 2025. *Qwen3-VL technical report*.
697 *Preprint*, arXiv:2511.21631.

698 Louisa Canepa, Sonit Singh, and Arcot Sowmya. 2023.
699 Visual question answering in the medical domain.
700 *arXiv preprint arXiv:2309.11080*.

701 Hao Chen, Wei Zhao, Yingli Li, Wenjun Li, Zhuoyi Li,
702 Ning Zhu, Tianyang Zhong, Yisong Wang, Youlan
703 Shang, Lei Guo, Junwei Han, Tianming Liu, Jun Liu,
704 and Tuo Zhang. 2025a. 3d-CT-GPT++: Enhancing
705 3d radiology report generation with direct preference
706 optimization and large vision-language models.

707 Liang Chen, Lei Li, Haozhe Zhao, Yifan Song, and
708 Vinci. 2025b. R1-v: Reinforcing super generaliza-
709 tion ability in vision-language models with less than
710 \$3. <https://github.com/Deep-Agent/R1-V>.

711 Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang
712 Zang, Zehui Chen, Haodong Duan, Jiaqi Wang,

713 Yu Qiao, Dahua Lin, and Feng Zhao. 2024a. *Are we
714 on the right way for evaluating large vision-language
715 models?* *Preprint*, arXiv:2403.20330.

716 Zhaorun Chen, Yichao Du, Zichen Wen, Yiyang Zhou,
717 Chenhang Cui, Zhenzhen Weng, Haoqin Tu, Chaoqi
718 Wang, Zhengwei Tong, Qinglan Huang, Canyu Chen,
719 Qinghao Ye, Zhihong Zhu, Yuqing Zhang, Jiawei
720 Zhou, Zhuokai Zhao, Rafael Rafailov, Chelsea Finn,
721 and Huaxiu Yao. 2024b. *Mj-bench: Is your multi-
722 modal reward model really a good judge for text-to-
723 image generation?* *Preprint*, arXiv:2407.04842.

724 Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye,
725 Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi
726 Hu, Jiapeng Luo, Zheng Ma, et al. 2024c. How far
727 are we to gpt-4v? closing the gap to commercial
728 multimodal models with open-source suites. *arXiv
729 preprint arXiv:2404.16821*.

730 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su,
731 Guo Chen, Sen Xing, Muyan Zhong, Qinglong
732 Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo,
733 Tong Lu, Yu Qiao, and Jifeng Dai. 2023. Internvl:
734 Scaling up vision foundation models and aligning
735 for generic visual-linguistic tasks. *arXiv preprint
736 arXiv:2312.14238*.

737 Kanzhi Cheng, Yantao Li, Fangzhi Xu, Jianbing Zhang,
738 Hao Zhou, and Yang Liu. 2024. Vision-language
739 models can self-improve reasoning via reflection.
740 *arXiv preprint arXiv:2411.00855*.

741 Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang
742 Tong, Saining Xie, Dale Schuurmans, Quoc V Le,
743 Sergey Levine, and Yi Ma. 2025. Sft memorizes,
744 rl generalizes: A comparative study of foundation
745 model post-training. *arXiv*.

746 Gheorghe Comanici, Eric Bieber, Mike Schaeckermann,
747 Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-
748 ccel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al.
749 2025. Gemini 2.5: Pushing the frontier with ad-
750 vanced reasoning, multimodality, long context, and
751 next generation agentic capabilities. *arXiv preprint
752 arXiv:2507.06261*.

753 Ganqu Cui, Lifan Yuan, Zefan Wang, Hanbin Wang,
754 Wendi Li, Bingxiang He, Yuchen Fan, Tianyu
755 Yu, Qixin Xu, Weize Chen, Jiarui Yuan, Huayu
756 Chen, Kaiyan Zhang, Xingtai Lv, Shuo Wang,
757 Yuan Yao, Xu Han, Hao Peng, Yu Cheng, Zhiyuan
758 Liu, Maosong Sun, Bowen Zhou, and Ning Ding.
759 2025. *Process reinforcement through implicit re-
760 wards*. *Preprint*, arXiv:2502.01456.

761 Wenliang Dai, Nayeon Lee, Boxin Wang, Zhuolin Yang,
762 Zihan Liu, Jon Barker, Tuomas Rintamaki, Moham-
763 mad Shoeybi, Bryan Catanzaro, and Wei Ping. 2024.
764 Nvlm: Open frontier-class multimodal llms. *arXiv
765 preprint arXiv:2409.11402*.

766 DeepSeek-AI. 2024. *Deepseek-v3 technical report*.
767 *Preprint*, arXiv:2412.19437.

768	Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. 2024. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. <i>arXiv preprint arXiv:2409.17146</i> .	Srishti Gureja, Lester James V. Miranda, Shayekh Bin Islam, Rishabh Maheshwary, Drishti Sharma, Gusti Winata, Nathan Lambert, Sebastian Ruder, Sara Hooker, and Marzieh Fadaee. 2024. <i>M-rewardbench: Evaluating reward models in multilingual settings</i> . Preprint, arXiv:2410.15522.	821
769			822
770			823
771			824
772			825
773			826
774	Zane Durante, Qiuyuan Huang, Naoki Wake, Ran Gong, Jae Sung Park, Bidipta Sarkar, Rohan Taori, Yusuke Noda, Demetri Terzopoulos, Yejin Choi, Katsushi Ikeuchi, Hoi Vo, Li Fei-Fei, and Jianfeng Gao. 2024. <i>Agent ai: Surveying the horizons of multimodal interaction</i> . Preprint, arXiv:2401.03568.	Lehan He, Zeren Chen, Zhelun Shi, Tianyu Yu, Jing Shao, and Lu Sheng. 2025. Systematic reward gap optimization for mitigating vlm hallucinations. <i>arXiv preprint arXiv:2411.17265</i> .	827
775			828
776			829
777			830
778			
779		Xuhao Hu, Dongrui Liu, Hao Li, Xuanjing Huang, and Jing Shao. 2024. Vlsbench: Unveiling visual leakage in multimodal safety. <i>arXiv preprint arXiv:2411.19939</i> .	831
780	SooHwan Eom, Jay Shim, Eunseop Yoon, Hee Suk Yoon, Hyeonmok Ko, Mark A. Hasegawa-Johnson, and Chang D. Yoo. 2025. SQuba: Speech mamba language model with querying-attention for efficient summarization.		832
781			833
782			834
783		Wen Huang, Hongbin Liu, Minxin Guo, and Neil Zhenqiang Gong. 2024. Visual hallucinations of multimodal large language models. <i>arXiv:2402.14683</i> .	835
784			836
			837
785	Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen, Meng Zhao, Yifan Zhang, Shaoqi Dong, Xiong Wang, Di Yin, Long Ma, et al. 2024a. Vita: Towards open-source interactive omni multimodal llm. <i>arXiv preprint arXiv:2408.05211</i> .	Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan, Zhonghao He, Lukas Vierling, Donghai Hong, Jiayi Zhou, Zhaowei Zhang, Fanzhi Zeng, Juntao Dai, Xuehai Pan, Kwan Yee Ng, Aidan O’Gara, Hua Xu, Brian Tse, Jie Fu, Stephen McAleer, Yaodong Yang, Yizhou Wang, Song-Chun Zhu, Yike Guo, and Wen Gao. 2025. <i>Ai alignment: A comprehensive survey</i> . Preprint, arXiv:2310.19852.	838
786			839
787			840
788			841
789			842
			843
			844
790	Chaoyou Fu, Haojia Lin, Xiong Wang, Yi-Fan Zhang, Yunhang Shen, Xiaoyu Liu, Yangze Li, Zuwei Long, Heting Gao, Ke Li, et al. 2025a. Vita-1.5: Towards gpt-4o level real-time vision and speech interaction. <i>arXiv preprint arXiv:2501.01957</i> .		845
791			846
792			
793			
794		Jiaming Ji, Jiayi Zhou, Hantao Lou, Boyuan Chen, Donghai Hong, Xuyao Wang, Wenqi Chen, Kaile Wang, Rui Pan, Jiahao Li, Mohan Wang, Josef Dai, Tianyi Qiu, Hua Xu, Dong Li, Weipeng Chen, Jun Song, Bo Zheng, and Yaodong Yang. 2024. <i>Align anything: Training all-modality models to follow instructions with language feedback</i> . Preprint, arXiv:2412.15838.	847
795	Jinlan Fu, Shenzhen Huangfu, Hao Fei, Yichong Huang, Xiaoyu Shen, Xipeng Qiu, and See-Kiong Ng. 2025b. Mcm-dpo: Multifaceted cross-modal direct preference optimization for alt-text generation. <i>arXiv preprint arXiv:2510.00647</i> .		848
796			849
797			850
798			851
799			852
			853
			854
800	Jinlan Fu, Shenzhen Huangfu, Hao Fei, Xiaoyu Shen, Bryan Hooi, Xipeng Qiu, and See-Kiong Ng. 2025c. <i>Chip: Cross-modal hierarchical direct preference optimization for multimodal llms</i> . Preprint, arXiv:2501.16629.	Hongrui Jia, Chaoya Jiang, Haiyang Xu, Wei Ye, Mengfan Dong, Ming Yan, Ji Zhang, Fei Huang, and Shikun Zhang. 2024. <i>Symdpo: Boosting in-context learning of large multimodal models with symbol demonstration direct preference optimization</i> . Preprint, arXiv:2411.11909.	855
801			856
802			857
803			858
804			859
			860
805	Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. 2024b. Blink: Multimodal large language models can see but not perceive. <i>arXiv preprint arXiv:2404.12390</i> .	Dongfu Jiang, Xuan He, Huaye Zeng, Cong Wei, Max Ku, Qian Liu, and Wenhui Chen. 2024. <i>Mantis: Interleaved multi-image instruction tuning</i> . Preprint, arXiv:2405.01483.	861
806			862
807			863
808			864
809			
810	Yuhan Fu, Ruobing Xie, Xingwu Sun, Zhanhui Kang, and Xirong Li. 2024c. <i>Mitigating hallucination in multimodal large language model via hallucination-targeted direct preference optimization</i> . Preprint, arXiv:2411.10436.	Songtao Jiang, Yan Zhang, Ruizhe Chen, Tianxiang Hu, Yeying Jin, Qinglin He, Yang Feng, Jian Wu, and Zuozhu Liu. 2025a. Modality-fair preference optimization for trustworthy mllm alignment. <i>arXiv preprint arXiv:2410.15334</i> .	865
811			866
812			867
813			868
814			869
815	Leo Gao, John Schulman, and Jacob Hilton. 2023. Scaling laws for reward model overoptimization. In <i>ICML</i> .	Songtao Jiang, Yan Zhang, Yeying Jin, Zhihang Tang, Yangyang Wu, Yang Feng, Jian Wu, and Zuozhu Liu. 2025b. Hscr: Hierarchical self-contrastive rewarding for aligning medical vision language models. <i>arXiv preprint arXiv:2506.00805</i> .	870
816			871
817			872
			873
			874
818	Anisha Gunjal, Jihan Yin, and Erhan Bas. 2024. <i>Detecting and preventing hallucinations in large vision language models</i> . Preprint, arXiv:2308.06394.	Kechen Jiao, Zhirui Fang, Jiahao Liu, Bei Li, Zhongjian Qiao, Yaxin Xu, Yifan Zhu, Xinyu Liu, Jingang	875
819			876
820			

877	Wang, and Xiu Li. 2025. InteractiveCOT: Aligning dynamic chain-of-thought planning for embodied decision-making.	930
878		931
879		932
880	Komal Kumar, Tajamul Ashraf, Omkar Thawakar, Rao Muhammad Anwer, Hisham Cholakkal, Mubarak Shah, Ming-Hsuan Yang, Phillip H. S. Torr, Fahad Shahbaz Khan, and Salman Khan. 2025. Llm post-training: A deep dive into reasoning large language models. <i>Preprint</i> , arXiv:2502.21321.	933
881		934
882		935
883		936
884		937
885		938
886	Kwai Keye Team, Biao Yang, Bin Wen, et al. 2025. Kwai Keye-VL technical report. <i>Preprint</i> , arXiv:2507.01949.	940
887		941
888		942
889	Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. 2024. Rewardbench: Evaluating reward models for language modeling. <i>Preprint</i> , arXiv:2403.13787.	943
890		944
891		945
892		946
893		947
894		948
895	Byung-Kwan Lee, Sangyun Chung, Chae Won Kim, Beomchan Park, and Yong Man Ro. 2024a. Phantom of latent for large language and vision models. <i>Preprint</i> , arXiv:2409.14713.	949
896		950
897		951
898		952
899	Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, and Sushant Prakash. 2024b. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback. <i>arXiv preprint arXiv:2309.00267</i> .	953
900		954
901		955
902		956
903		957
904		958
905	Kang-il Lee, Minbeom Kim, Seunghyun Yoon, Min-sung Kim, Dongryeol Lee, Hyukhun Koh, and Kyomin Jung. 2024c. Vblind-bench: Measuring language priors in large vision-language models. <i>arXiv:2406.08702</i> .	959
906		960
907		961
908		962
909		963
910	Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. 2024a. Llava-next: Stronger llms supercharge multimodal capabilities in the wild.	964
911		965
912		966
913		967
914	Dongping Li, Tielong Cai, Tianci Tang, Wenhao Chai, Katherine Rose Driggs-Campbell, and Gaoang Wang. 2025a. Emmoe: A comprehensive benchmark for embodied mobile manipulation in open environments. <i>Preprint</i> , arXiv:2503.08604.	968
915		969
916		970
917		971
918		972
919	Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. 2024b. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. <i>Preprint</i> , arXiv:2407.07895.	973
920		974
921		975
922		976
923		977
924	Lei Li, Yuancheng Wei, Zhihui Xie, Xuqing Yang, Yifan Song, Peiyi Wang, Chenxin An, Tianyu Liu, Sujian Li, Bill Yuchen Lin, et al. 2024c. Vrewardbench: A challenging benchmark for vision-language generative reward models. <i>arXiv preprint arXiv:2411.17451</i> .	978
925		979
926		980
927		981
928		982
929		983
		984
	Lei Li, Zhihui Xie, Mukai Li, Shunian Chen, Peiyi Wang, Liang Chen, Yazheng Yang, Benyou Wang, and Lingpeng Kong. 2023a. Silkie: Preference distillation for large visual language models. <i>arXiv preprint arXiv:2312.10665</i> .	930
		931
		932
		933
	Xiaoqi Li, Mingxu Zhang, Yiran Geng, Haoran Geng, Yuxing Long, Yan Shen, Renrui Zhang, Jiaming Liu, and Hao Dong. 2023b. Manipllm: Embodied multimodal large language model for object-centric robotic manipulation. <i>Preprint</i> , arXiv:2312.16217.	935
		936
		937
		938
		939
	Xudong Li, Mengdan Zhang, Peixian Chen, Xiawu Zheng, Yan Zhang, Jingyuan Zheng, Yunhang Shen, Ke Li, Chaoyou Fu, Xing Sun, and Rongrong Ji. 2025b. Zooming from context to cue: Hierarchical preference optimization for multi-image mllms. <i>arXiv preprint arXiv:2505.22396</i> .	940
		941
		942
		943
		944
		945
	Yongqi Li, Lu Yang, Jian Wang, Runyang You, Wenjie Li, and Liqiang Nie. 2025c. Towards harmless multimodal assistants with blind preference optimization. <i>arXiv preprint arXiv:2503.14189</i> .	946
		947
		948
		949
	Zongxia Li, Xiyang Wu, Hongyang Du, Fuxiao Liu, Huy Nghiem, and Guangyao Shi. 2025d. A survey of state of the art large vision language models: Alignment, benchmark, evaluations and challenges. <i>Preprint</i> , arXiv:2501.02189.	950
		951
		952
		953
		954
	Jiafeng Liang, Shixin Jiang, Xuan Dong, Ning Wang, Zheng Chu, Hui Su, Jinlan Fu, Ming Liu, See-Kiong Ng, and Bing Qin. 2025a. Investigating and enhancing the robustness of large multimodal models against temporal inconsistency. <i>arXiv preprint arXiv:2505.14405</i> .	955
		956
		957
		958
		959
		960
	Wenlong Liang, Rui Zhou, Yang Ma, Bing Zhang, Songlin Li, Yijia Liao, and Ping Kuang. 2025b. Large model empowered embodied ai: A survey on decision-making and embodied learning. <i>Preprint</i> , arXiv:2508.10399.	961
		962
		963
		964
		965
	Chaohu Liu, Gui Tianyi, Yu Liu, and Linli Xu. 2024a. AdPO: Enhancing the adversarial robustness of large vision-language models with preference optimization.	966
		967
		968
		969
	Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. 2023a. Mitigating hallucination in large multi-modal models via robust instruction tuning. In <i>ICLR</i> .	970
		971
		972
		973
	Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. 2023b. Mitigating hallucination in large multi-modal models via robust instruction tuning. In <i>The Twelfth International Conference on Learning Representations</i> .	974
		975
		976
		977
		978
	Jiaming Liu, Chenxuan Li, Guanqun Wang, Lily Lee, Kaichen Zhou, Sixiang Chen, Chuyan Xiong, Jiaxin Ge, Renrui Zhang, and Shanghang Zhang. 2024b. Self-corrected multimodal large language model for end-to-end robot manipulation. <i>Preprint</i> , arXiv:2405.17418.	979
		980
		981
		982
		983
		984

985	Ruyang Liu, Chen Li, Haoran Tang, Yixiao Ge, Ying Shan, Haibo Lu, and Jiankun Yang. 2025a. PPLLaVA: Varied video sequence understanding with prompt guidance.	Timothy Ossowski, Jixuan Chen, Danyal Maqbool, Zefan Cai, Tyler Bradshaw, and Junjie Hu. 2025. Comma: A communicative multimodal multi-agent benchmark. <i>Preprint</i> , arXiv:2410.07553.	1037 1038 1039 1040
989	Wenqi Liu, Xuemeng Song, Jiayi Li, Yinwei Wei, Na Zheng, Jianhua Yin, and Liqiang Nie. 2025b. Mitigating hallucination through theory-consistent symmetric multimodal preference optimization. <i>arXiv preprint arXiv:2506.11712</i> .	Yassine Ouali, Adrian Bulat, Brais Martinez, and Georgios Tzimiropoulos. 2024. Clip-dpo: Vision-language models as a source of preference for fixing hallucinations in vlms. <i>Preprint</i> , arXiv:2408.10433.	1041 1042 1043 1044
994	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. 2023c. Mmbench: Is your multi-modal model an all-around player? <i>arXiv preprint arXiv:2307.06281</i> .	Piotr Padlewski, Max Bain, Matthew Henderson, Zhongkai Zhu, Nishant Relan, Hai Pham, Donovan Ong, Kaloyan Aleksiev, Aitor Ormazabal, Samuel Phua, Ethan Yeo, Eugenie Lamprecht, Qi Liu, Yuqi Wang, Eric Chen, Deyu Fu, Lei Li, Che Zheng, Cyprien de Masson d’Autume, Dani Yogatama, Mikel Artetxe, and Yi Tay. 2024. Vibe-eval: A hard evaluation suite for measuring progress of multimodal language models. <i>Preprint</i> , arXiv:2405.02287.	1045 1046 1047 1048 1049 1050 1051 1052 1053
999	Ziyu Liu, Yuhang Zang, Xiaoyi Dong, Pan Zhang, Yuhang Cao, Haodong Duan, Conghui He, Yuanjun Xiong, Dahua Lin, and Jiaqi Wang. 2024c. Mia-dpo: Multi-image augmented direct preference optimization for large vision-language models. <i>Preprint</i> , arXiv:2410.17637.	Han Qiu, Peng Gao, Lewei Lu, Xiaoqin Zhang, Ling Shao, and Shijian Lu. 2025. Spatial preference rewarding for mllms spatial understanding. <i>arXiv preprint arXiv:2510.14374</i> .	1054 1055 1056 1057
1005	Jinda Lu, Junkang Wu, Jinghan Li, Xiaojun Jia, Shuo Wang, YiFan Zhang, Junfeng Fang, Xiang Wang, and Xiangnan He. 2025. Dama: Data- and model-aware alignment of multi-modal llms. <i>Preprint</i> , arXiv:2502.01943.	Rafael Rafailov, Yaswanth Chittipedu, Ryan Park, Harshit Sikchi, Joey Hejna, W. Bradley Knox, Chelsea Finn, and Scott Niekum. 2024. Scaling laws for reward model overoptimization in direct alignment algorithms. <i>CoRR</i> .	1058 1059 1060 1061 1062
1010	Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. <i>arXiv:2310.02255</i> .	Yifan Shen, Yuanzhe Liu, Jingyuan Zhu, Xu Cao, Xiaofeng Zhang, Yixiao He, Wenming Ye, James Matthew Rehg, and Ismini Lourentzou. 2025. Fine-grained preference optimization improves spatial reasoning in vlms. <i>arXiv preprint arXiv:2506.21656</i> .	1063 1064 1065 1066 1067 1068
1016	Tiange Luo, Ang Cao, Gunhee Lee, Justin Johnson, and Honglak Lee. 2025. vVLM: Exploring visual reasoning in VLMs against language priors.	Dong Shu, Haiyan Zhao, Jingyu Hu, Weiru Liu, Ali Payani, Lu Cheng, and Mengnan Du. 2025. Large vision-language model alignment and misalignment: A survey through the lens of explainability. In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 1713–1735, Suzhou, China. Association for Computational Linguistics.	1069 1070 1071 1072 1073 1074 1075
1019	Feipeng Ma, Yizhou Zhou, Yueyi Zhang, Siying Wu, Zheyu Zhang, Zilong He, Fengyun Rao, and Xiaoyan Sun. 2024a. Task navigator: Decomposing complex tasks for multimodal large language models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> .	Zihan Song, Xin Wang, Zi Qian, Hong Chen, Longtao Huang, Hui Xue, and Wenwu Zhu. 2025. Modularized self-reflected video reasoner for multimodal LLM with application to video question answering. In <i>Forty-second International Conference on Machine Learning</i> .	1076 1077 1078 1079 1080 1081
1025	Xiaojian Ma, Silong Yong, Zilong Zheng, Qing Li, Yitao Liang, Song-Chun Zhu, and Siyuan Huang. 2023. Sqa3d: Situated question answering in 3d scenes. In <i>ICLR</i> .	Guangzhi Sun, Yudong Yang, Jimin Zhuang, Changli Tang, Yixuan Li, Wei Li, Zejun MA, and Chao Zhang. 2025. video-salmonn-o1: Reasoning-enhanced audio-visual large language model. <i>arXiv preprint arXiv:2502.11775</i> .	1082 1083 1084 1085 1086
1029	Yueen Ma, Zixing Song, Yuzheng Zhuang, Jianye Hao, and Irwin King. 2024b. A survey on vision-language-action models for embodied ai. <i>Preprint</i> , arXiv:2405.14093.	Guangzhi Sun, Wenyi Yu, Changli Tang, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, Yuxuan Wang, and Chao Zhang. 2024. video-salmonn: Speech-enhanced audio-visual large language models. In <i>International Conference on Machine Learning</i> , pages 47198–47217. PMLR.	1087 1088 1089 1090 1091 1092
1033	Yu Meng, Mengzhou Xia, and Danqi Chen. 2024. Simpo: Simple preference optimization with a reference-free reward. <i>CoRR</i> .		
1036	OpenAI. 2024. Introducing openai o1-preview .		

1093	Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, Kurt Keutzer, and Trevor Darrell. 2023a. Aligning large multimodal models with factually augmented rlhf . <i>Preprint</i> , arXiv:2309.14525.		
1094		Zhu, Lewei Lu, Yu Qiao, and Jifeng Dai. 2024d. Enhancing the reasoning ability of multimodal large language models via mixed preference optimization . <i>Preprint</i> , arXiv:2411.10442.	1147
1095			1148
1096			1149
1097			1150
1098		Weiyun Wang et al. 2025c. InternVL3.5: Advancing open-source multimodal models in versatility, reasoning, and efficiency . <i>Preprint</i> , arXiv:2508.18265.	1151
1099	Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. 2023b. Aligning large multimodal models with factually augmented rlhf . <i>arXiv:2309.14525</i> .		1152
1100			1153
1101		Xinming Wang, Jian Xu, Aslan H Feng, Yi Chen, Haiyang Guo, Fei Zhu, Yuanqi Shao, Minsi Ren, Hongzhu Yi, Sheng Lian, et al. 2025d. The hitchhiker’s guide to autonomous research: A survey of scientific agents. <i>TechRxiv.August 07, 2025</i> . DOI:10.36227/techrxiv175459840.02185500/V1.	1154
1102			1155
1103			1156
1104	Changli Tang, Yixuan Li, Yudong Yang, Jimin Zhuang, Guangzhi Sun, Wei Li, Zejun MA, and Chao Zhang. 2025. Enhancing multimodal LLM for detailed and accurate video captioning using multi-round preference optimization .		1157
1105			1158
1106			1159
1107		Xinming Wang, Jian Xu, Bin Yu, Sheng Lian, Hongzhu Yi, Yi Chen, Yingjian Zhu, Boran Wang, Hongming Yang, Han Hu, et al. 2025e. Mr-align: Meta-reasoning informed factuality alignment for large reasoning models . <i>arXiv preprint arXiv:2510.24794</i> .	1160
1108			1161
1109	Changli Tang, Yixuan Li, Yudong Yang, Jimin Zhuang, Guangzhi Sun, Wei Li, Zujun Ma, and Chao Zhang. 2024. Enhancing multimodal llm for detailed and accurate video captioning using multi-round preference optimization . <i>arXiv preprint arXiv:2410.06682</i> .		1162
1110			1163
1111			1164
1112		Xiyao Wang, Jiu Hai Chen, Zhaoyang Wang, Yuhang Zhou, Yiyang Zhou, Huaxiu Yao, Tianyi Zhou, Tom Goldstein, Parminder Bhatia, Furong Huang, and Cao Xiao. 2024e. Enhancing visual-language modality alignment in large vision language models via self-improvement . <i>Preprint</i> , arXiv:2405.15973.	1165
1113			1166
1114			1167
1115	Gemini Team. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context . <i>Preprint</i> , arXiv:2403.05530.		1168
1116			1169
1117			1170
1118	Shubham Toshniwal, Ivan Moshkov, Sean Narenthiran, Daria Gitman, Fei Jia, and Igor Gitman. 2024. Openmathinstruct-1: A 1.8 million math instruction tuning dataset . <i>Preprint</i> , arXiv:2402.10176.		1171
1119			1172
1120			1173
1121			1174
1122			1175
1123	Fei Wang, Wenxuan Zhou, James Y. Huang, Nan Xu, Sheng Zhang, Hoifung Poon, and Muhao Chen. 2024a. mdpo: Conditional preference optimization for multimodal large language models . <i>Preprint</i> , arXiv:2406.11839.		1176
1124			1177
1125			1178
1126			1179
1127			1180
1128	Han Wang, Gang Wang, and Huan Zhang. 2025a. Steering away from harm: An adaptive approach to defending vision language model against jailbreaks . <i>arXiv preprint arXiv:2411.16721</i> .		1181
1129			1182
1130		Colin White, Samuel Dooley, Manley Roberts, Arka Pal, Ben Feuer, Siddhartha Jain, Ravid Shwartz-Ziv, Neel Jain, Khalid Saifullah, Siddhartha Naidu, Chinmay Hegde, Yann LeCun, Tom Goldstein, Willie Neiswanger, and Micah Goldblum. 2024. Livebench: A challenging, contamination-free llm benchmark . <i>Preprint</i> , arXiv:2406.19314.	1183
1131			1184
1132			1185
1133			1186
1134	Junyong Wang, Haiyang Xu, Jiabo Ye, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang. 2024b. Mobile-agent: Autonomous multi-modal mobile device agent with visual perception . <i>Preprint</i> , arXiv:2401.16158.		1187
1135			1188
1136			1189
1137			1190
1138			1191
1139			1192
1140			1193
1141			1194
1142			1195
1143			1196
1144			1197
1145			1198
1146			1199
			1200
			1201
			1202
			1203
			1204
			1205
			1206
			1207
			1208
			1209
			1210
			1211
			1212
			1213
			1214
			1215
			1216
			1217
			1218
			1219
			1220
			1221
			1222
			1223
			1224
			1225
			1226
			1227
			1228
			1229
			1230
			1231
			1232
			1233
			1234
			1235
			1236
			1237
			1238
			1239
			1240
			1241
			1242
			1243
			1244
			1245
			1246
			1247
			1248
			1249
			1250
			1251
			1252
			1253
			1254
			1255
			1256
			1257
			1258
			1259
			1260
			1261
			1262
			1263
			1264
			1265
			1266
			1267
			1268
			1269
			1270
			1271
			1272
			1273
			1274
			1275
			1276
			1277
			1278
			1279
			1280
			1281
			1282
			1283
			1284
			1285
			1286
			1287
			1288
			1289
			1290
			1291
			1292
			1293
			1294
			1295
			1296
			1297
			1298
			1299
			1300
			1301
			1302

1203	Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui.	trustworthy mllms via behavior alignment from fine-	1258
1204	2023. The rise and potential of large language model	grained correctional human feedback. In <i>Proceed-</i>	1259
1205	based agents: A survey . <i>Preprint</i> , arXiv:2309.07864.	<i>ings of the IEEE/CVF Conference on Computer Vi-</i>	1260
		<i>sion and Pattern Recognition</i> .	1261
1206	Tianyi Xiong, Xiyao Wang, Dong Guo, Qinghao Ye,	Tianyu Yu, Haoye Zhang, Yuan Yao, Yunkai Dang,	1262
1207	Haoqi Fan, Quanquan Gu, Heng Huang, and Chun-	Da Chen, Xiaoman Lu, Ganqu Cui, Taiwen He,	1263
1208	yuan Li. 2024. Llava-critic: Learning to evaluate	Zhiyuan Liu, Tat-Seng Chua, et al. 2024b. Rlaif-	1264
1209	multimodal models . <i>Preprint</i> , arXiv:2410.02712.	v: Aligning mllms through open-source ai feedback	1265
		for super gpt-4v trustworthiness. <i>arXiv preprint</i>	1266
1210	Chejian Xu, Jiawei Zhang, Zhaorun Chen, Chulin Xie,	<i>arXiv:2405.17220</i> .	1267
1211	Mintong Kang, Yujin Potter, Zhun Wang, Zhuowen		
1212	Yuan, Alexander Xiong, Zidi Xiong, Chenhui Zhang,	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng,	1268
1213	Lingzhi Yuan, Yi Zeng, Peiyang Xu, Chengquan	Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang,	1269
1214	Guo, Andy Zhou, Jeffrey Ziwei Tan, Xuandong	Weiming Ren, Yuxuan Sun, et al. 2024. Mmmu: A	1270
1215	Zhao, Francesco Pinto, Zhen Xiang, Yu Gai, Zi-	massive multi-discipline multimodal understanding	1271
1216	nan Lin, Dan Hendrycks, Bo Li, and Dawn Song.	and reasoning benchmark for expert agi. In <i>CVPR</i> .	1272
1217	2025. Mmdt: Decoding the trustworthiness and		
1218	safety of multimodal foundation models. <i>arXiv</i>	Fatemeh Pesaran Zadeh, Yoojin Oh, and Gunhee	1273
1219	<i>preprint arXiv:2503.14827</i> .	Kim. 2025. Lpoi: Listwise preference optimiza-	1274
		tion for vision language models. <i>arXiv preprint</i>	1275
1220	Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan,	<i>arXiv:2505.21061</i> .	1276
1221	Senthil Purushwalkam, Honglu Zhou, Viraj Prabhu,		
1222	Yutong Dai, Michael S Ryoo, Shrikant Kendre, Jieyu	Yuexiang Zhai, Hao Bai, Zipeng Lin, Jiayi Pan, Sheng-	1277
1223	Zhang, Can Qin, Shu Zhang, Chia-Chih Chen, Ning	bang Tong, Yifei Zhou, Alane Suhr, Saining Xie,	1278
1224	Yu, Juntao Tan, Tulika Manoj Awalgaoonkar, Shelby	Yann LeCun, Yi Ma, and Sergey Levine. 2024. Fine-	1279
1225	Heinecke, Huan Wang, Yejin Choi, Ludwig Schmidt,	tuning large vision-language models as decision-	1280
1226	Zeyuan Chen, Silvio Savarese, Juan Carlos Niebles,	making agents via reinforcement learning . <i>Preprint</i> ,	1281
1227	Caiming Xiong, and Ran Xu. 2024. xgen-mm (blip-	arXiv:2405.10292.	1282
1228	3): A family of open large multimodal models .		
1229	<i>Preprint</i> , arXiv:2408.08872.	Menghao Zhang, Huazheng Wang, Pengfei Ren,	1283
		Kangheng Lin, Qi Qi, Haifeng Sun, Zirui Zhuang,	1284
1230	Bei Yan, Jie Zhang, Zheng Yuan, Shiguang Shan, and	Lei Zhang, Jianxin Liao, and Jingyu Wang. 2025a.	1285
1231	Xilin Chen. 2024. Evaluating the quality of halluci-	Do LVLMs truly understand video anomalies? re-	1286
1232	nation benchmarks for large vision-language models.	vealing hallucination via co-occurrence patterns . In	1287
1233	<i>arXiv:2406.17115</i> .	<i>The Thirty-ninth Annual Conference on Neural Infor-</i>	1288
		<i>mation Processing Systems</i> .	1289
1234	An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao,	Renrui Zhang, Xinyu Wei, Dongzhi Jiang, Ziyu Guo,	1290
1235	Bowen Yu, Chengpeng Li, Dayiheng Liu, Jian-	Shicheng Li, Yichi Zhang, Chengzhuo Tong, Jiaming	1291
1236	hong Tu, Jingren Zhou, Junyang Lin, Keming Lu,	Liu, Aojun Zhou, Bin Wei, Shanghang Zhang, Peng	1292
1237	Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang	Gao, Chunyuan Li, and Hongsheng Li. 2024a. Mavis:	1293
1238	Ren, and Zhenru Zhang. 2024a. Qwen2.5-math tech-	Mathematical visual instruction tuning with an auto-	1294
1239	nical report: Toward mathematical expert model via	matic data engine . <i>Preprint</i> , arXiv:2407.08739.	1295
1240	self-improvement . <i>Preprint</i> , arXiv:2409.12122.		
		Ruohong Zhang, Bowen Zhang, Yanghao Li, Haotian	1296
1241	Yulong Yang, Xinshan Yang, Shuaidong Li, Chenhao	Zhang, Zhiqing Sun, Zhe Gan, Yinfei Yang, Ruom-	1297
1242	Lin, Zhengyu Zhao, Chao Shen, and Tianwei Zhang.	ing Pang, and Yiming Yang. 2024b. Improve vision	1298
1243	2024b. Security matrix for multimodal agents on	language model chain-of-thought reasoning. <i>arXiv</i>	1299
1244	mobile devices: A systematic and proof of concept	<i>preprint arXiv:2410.16198</i> .	1300
1245	study . <i>Preprint</i> , arXiv:2407.09295.		
		Yi-Fan Zhang, Tao Yu, Haochen Tian, Chaoyou Fu,	1301
1246	Kaining Ying, Fanqing Meng, Jin Wang, Zhiqian Li,	Peiyan Li, Jianshu Zeng, Wulin Xie, Yang Shi,	1302
1247	Han Lin, Yue Yang, Hao Zhang, Wenbo Zhang, Yuqi	Huanyu Zhang, Junkang Wu, Xue Wang, Yibo Hu,	1303
1248	Lin, Shuo Liu, Jiayi Lei, Quanfeng Lu, Runjian Chen,	Bin Wen, Fan Yang, Zhang Zhang, Tingting Gao,	1304
1249	Peng Xu, Renrui Zhang, Haozhe Zhang, Peng Gao,	Di Zhang, Liang Wang, Rong Jin, and Tieniu Tan.	1305
1250	Yali Wang, Yu Qiao, Ping Luo, Kaipeng Zhang, and	2025b. Mm-rlhf: The next step forward in multi-	1306
1251	Wenqi Shao. 2024. Mmt-bench: A comprehensive	modal llm alignment. <i>arXiv</i> .	1307
1252	multimodal benchmark for evaluating large vision-		
1253	language models towards multitask agi . <i>Preprint</i> ,	Yi-Fan Zhang, Weichen Yu, Qingsong Wen, Xue Wang,	1308
1254	arXiv:2404.16006.	Zhang Zhang, Liang Wang, Rong Jin, and Tieniu	1309
		Tan. 2024c. Debiasing large visual language models.	1310
		<i>arXiv preprint arXiv:2403.05262</i> .	1311
1255	Tianyu Yu, Yuan Yao, Haoye Zhang, Taiwen He, Yifeng	Yi-Fan Zhang, Huanyu Zhang, Haochen Tian, Chaoyou	1312
1256	Han, Ganqu Cui, Jinyi Hu, Zhiyuan Liu, Hai-Tao	Fu, Shuangqing Zhang, Junfei Wu, Feng Li,	1313
1257	Zheng, Maosong Sun, et al. 2024a. Rllhf-v: Towards		

1314 Kun Wang, Qingsong Wen, Zhang Zhang, et al.
1315 2024d. Mme-realworld: Could your multimodal
1316 llm challenge high-resolution real-world scenarios
1317 that are difficult for humans? *arXiv preprint*
1318 *arXiv:2408.13257*.

1319 Yichi Zhang, Yao Huang, Yitong Sun, Chang Liu, Zhe
1320 Zhao, Zhengwei Fang, Yifan Wang, Huanran Chen,
1321 Xiao Yang, Xingxing Wei, et al. Multitrust: A com-
1322 prehensive benchmark towards trustworthy multi-
1323 modal large language models. In *The Thirty-eight*
1324 *Conference on Neural Information Processing Sys-*
1325 *tems Datasets and Benchmarks Track*.

1326 Zhihan Zhang, Yixin Cao, and Lizi Liao. 2025c.
1327 Boosting chart-to-code generation in mllm via
1328 dual preference-guided refinement. *arXiv preprint*
1329 *arXiv:2504.02906*.

1330 Zijian Zhang, Kaiyuan Zheng, Zhaorun Chen, Joel Jang,
1331 Yi Li, Siwei Han, Chaoqi Wang, Mingyu Ding, Di-
1332 eter Fox, and Huaxiu Yao. 2025d. [Grape: Generaliz-](#)
1333 [ing robot policy via preference alignment](#). *Preprint*,
1334 *arXiv:2411.19309*.

1335 Zhiyuan Zhao, Bin Wang, Linke Ouyang, Xiaoyi Dong,
1336 Jiaqi Wang, and Conghui He. 2024. [Beyond hallu-](#)
1337 [cinations: Enhancing lvlms through hallucination-](#)
1338 [aware direct preference optimization](#). *Preprint*,
1339 *arXiv:2311.16839*.

1340 Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao
1341 Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu,
1342 Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis,
1343 Luke Zettlemoyer, and Omer Levy. 2023. [Lima: Less](#)
1344 [is more for alignment](#). *Preprint*, *arXiv:2305.11206*.

1345 Ziyin Zhou, Yunpeng Luo, Yuanchen Wu, Ke Sun, Jiayi
1346 Ji, Ke Yan, Shouhong Ding, Xiaoshuai Sun, Yun-
1347 sheng Wu, and Rongrong Ji. 2025. [Aigi-holmes:](#)
1348 [Towards explainable and generalizable ai-generated](#)
1349 [image detection via multimodal large language mod-](#)
1350 [els](#). *arXiv preprint arXiv:2507.02664*.

1351 Kangyu Zhu, Peng Xia, Yun Li, Hongtu Zhu, Sheng
1352 Wang, and Huaxiu Yao. 2025a. [Mmedpo: Aligning](#)
1353 [medical vision-language models with clinical-aware](#)
1354 [multimodal preference optimization](#). *arXiv preprint*
1355 *arXiv:2412.06141*.

1356 Lanyun Zhu, Deyi Ji, Tianrun Chen, Haiyang Wu,
1357 De Wen Soh, and Jun Liu. 2025b. [CPCF: A cross-](#)
1358 [prompt contrastive framework for referring multi-](#)
1359 [modal large language models](#). In *Forty-second Inter-*
1360 *national Conference on Machine Learning*.

1361 Yongshuo Zong, Ondrej Bohdal, Tingyang Yu, Yongxin
1362 Yang, and Timothy Hospedales. 2024. [Safety fine-](#)
1363 [tuning at \(almost\) no cost: A baseline for vision large](#)
1364 [language models](#). *arXiv preprint arXiv:2402.02207*.

Appendix

A Preliminary of Optimization Paradigm 16

A.1 RLHF 16

A.2 DPO 16

B Background of MLLM Alignment 16

B.1 Pre-Training 16

B.2 Instruction Tuning 17

B.3 Alignment with Human Preference 17

C Positioning Relative to Prior Works 17

D Future Directions 17

E Leveraging Visual Information for Alignment 18

A Preliminary of Optimization Paradigm

A.1 RLHF

RLHF is a technique that combines RL with human feedback to align AI models (such as LLMs) with human preferences or values. It is widely used in alignment tasks and has played a crucial role in fine-tuning generative models (e.g., ChatGPT). RLHF typically consists of three stages: pretrained model, reward modeling, and RL fine-tuning.

In the field of MLLM alignment, RLHF-based methods are relatively rare, with Fact-RLHF being a representative example. Fact-RLHF undergoes three training phases (supervised fine-tuning, human preference collection & preference modeling, and factually-augmented RLHF) to align with human preferences. In Table 1, we present the loss function of the Fact-RLHF algorithm and compare it with other alignment algorithms. The training objective of RLHF is as follows:

$$\mathcal{L}_{\text{RLHF}} = -\mathbf{E}_{(\mathcal{I}, x) \in \mathcal{D}, y \sim \pi_{\phi}(y|\mathcal{I}, x)} [r_{\theta}(\mathcal{I}, x, y) - \beta \cdot \mathbb{D}_{KL}(\pi_{\phi}(y|\mathcal{I}, x) \parallel \pi^{\text{INIT}}(y|\mathcal{I}, x))], \quad (1)$$

where \mathcal{L} denotes the loss function, \mathcal{D} represents the dataset, \mathcal{I} represents the image, x represents the question, y represents the response, \mathbb{D}_{KL} represents the KL divergence, π_{ϕ} represents the trained MLLM policy, π^{INIT} represents the fixed initial policy model, r_{θ} represents the reward model, and β represents a hyperparameter.

A.2 DPO

DPO is an emerging model alignment approach designed to replace traditional RLHF by directly

optimizing model policies using preference data, eliminating the need for explicit reward model training or complex reinforcement learning algorithms. Its core idea is to transform human preferences into optimization objectives for probabilistic models, thereby simplifying the alignment process and improving efficiency.

Since DPO does not require an explicit reward model, it reduces the number of hyperparameters, making it widely adopted in the MLLM alignment field. In Table 1, we present the DPO loss function and all its variants in MLLM alignment, allowing readers to quickly and intuitively understand the differences and connections between different algorithms. The DPO training objective is defined as follows:

$$\mathcal{L}_{\text{dpo}} = -\mathbf{E}_{(\mathcal{I}, x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(\beta \log \frac{\pi_{\theta}(y_w|\mathcal{I}, x)}{\pi_{\text{ref}}(y_w|\mathcal{I}, x)}) - \beta \log \frac{\pi_{\theta}(y_l|\mathcal{I}, x)}{\pi_{\text{ref}}(y_l|\mathcal{I}, x)}), \quad (2)$$

where \mathcal{L} denotes the loss function, \mathcal{D} represents the dataset, \mathcal{I} represents the image, x represents the question, y_w, y_l represents the chosen and rejected responses respectively, \mathbb{D}_{KL} represents the KL divergence, π_{ϕ} represents the trained MLLM policy, π_{ref} represents the reference MLLM policy, r_{θ} represents the reward model, and β represents a hyperparameter, σ represents the sigmoid function.

B Background of MLLM Alignment

In this section, we will provide a brief explanation of the complete training process for MLLMs, which consists mainly of three phases: pre-training, instruction tuning, and alignment with human preference.

B.1 Pre-Training

The pre-training phase of MLLMs primarily aims to align the feature spaces of different modalities with that of the language model. The data used in this phase is typically simple caption data. For instance, image-caption pairs are commonly used for image/video understanding MLLMs (Bai et al., 2023; Chen et al., 2023), while speech data and transcriptions are used for speech understanding MLLMs (Fu et al., 2024a, 2025a). Through this pre-training phase, the model learns to understand inputs from various modalities.

1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460

1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480

1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498

B.2 Instruction Tuning

Building on the pre-training phase, the SFT phase aims to teach the model how to interact with humans by focusing on understanding questions and providing responses in a specified format, i.e., instruction-following ability. The data used in this phase is typically high-quality and diverse dialogue data. For example, in the commonly seen visual question answering (VQA) task, given an image and its corresponding instruction, the trained model will provide the correct answer for the task.

B.3 Alignment with Human Preference

Previous works have shown that SFT tends to make the model memorize training data and try to generalize across diverse scenarios (Chu et al., 2025). The alignment phase, typically involving reinforcement learning (RL) strategies, is crucial for generalizing to unseen domains. However, most multimodal models neglect this step (Wang et al., 2024c; Deitke et al., 2024; Chen et al., 2024c; Dai et al., 2024; Agrawal et al., 2024). The goals of the alignment stage are broad, such as reducing hallucinations (Zhang et al., 2024c; Lu et al., 2025), enhancing conversational abilities (Xiong et al., 2024), improving safety (Zong et al., 2024), strengthening the reasoning abilities (Wang et al., 2024d), improving capabilities for long-reasoning tasks like DeepSeek-R1 (Chen et al., 2025b), and overall MLLM performance (Zhang et al., 2025b). This phase usually uses pair data that incorporates human preference.

C Positioning Relative to Prior Works

Here, we reiterate the systematic and innovative aspects of our work compared to other surveys. First, previous alignment review papers have primarily focused on AI (Ji et al., 2025) or LLM (Kumar et al., 2025) alignment, with no existing work specifically addressing MLLM alignment. Our paper is the first survey dedicated to MLLM alignment. Among related works, (Kumar et al., 2025) is the most relevant, but it provides only a broad overview of LLM post-training methods without delving into their specific characteristics, making it more of a general summary. In contrast, we offer a detailed analysis of each method, enabling readers to intuitively understand their distinctions and quickly grasp the fundamentals of MLLM alignment. While MLLMs demonstrate immense potential in handling complex tasks involving visual,

auditory, and textual data, state-of-the-art MLLMs are rarely rigorously aligned with human preferences. Our paper aims to provide researchers with a comprehensive and systematic guide to entering the field of MLLM alignment.

Second, we systematically organize MLLM alignment algorithms based on their application scenarios, data construction, and evaluation methods, addressing key questions such as how algorithms are applied, how datasets are built, and how methods are evaluated. This forms a clear and structured framework for MLLM alignment:

(1) **Algorithms:** In Table 1, we establish a unified notation system to formally describe different algorithms, allowing readers to grasp their differences and connections quickly.

(2) **Datasets:** In Table 2, we summarize existing open-source datasets in terms of size, categories, response-generating models, data sources, and annotation models, providing readers with a clear and concise overview of MLLM alignment-specific datasets.

(3) **Technical Pipeline:** In Fig 1 and Fig 3, we summarize the end-to-end technical pipeline of MLLM alignment, covering benchmark-driven evaluation and preference construction, and highlight the key design choices and method characteristics at each stage to provide a practical reference for researchers.

Third, we thoroughly discuss emerging research directions and open challenges in MLLM alignment, such as visual challenges, comprehensive evaluation, full-modality alignment, MLLM reasoning, insights from LLM alignment, and MLLM as agents. These are critical issues uniquely faced by MLLM alignment today, and our paper is the first to systematically propose them in Appendix D. By doing so, we aim to inspire potential research ideas and foster growth within the alignment community.

D Future Directions

Despite rapid progress in multimodal post-training, several open challenges remain in scaling alignment beyond image–text, transferring reasoning-oriented RL or data advances to MLLMs, mitigating overoptimization and reward hacking, and enabling robust, secure agentic behaviors under multimodal interaction. Given these circumstances, we outline key directions spanning full-modality alignment, reasoning-centric optimization, insights

1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538

1539
1540
1541
1542
1543
1544
1545
1546
1547
1548

from LLM alignment, and MLLM-based agents.

Full-modality alignment Align-anything(Ji et al., 2024) pioneers full-modality alignment through the multimodal dataset "align-anything-200k", which spans text, images, audio, and video. This study demonstrates the complementary effects between different modalities. However, their work is still in its early stages. The dataset for each modality is relatively small, limiting its ability to cover a wide range of tasks. Additionally, the proposed algorithm is only a preliminary improvement on the DPO method, and it does not fully exploit the unique structural information inherent in each modality. Moving forward, the design of alignment algorithms beyond image/text domains, particularly for other modalities, to enhance multimodal model capabilities, will be a key trend.

MLLM reasoning Recent advancements in reasoning LLMs, such as OpenAI’s O1 and DeepSeek-R1, highlight the importance of RL algorithms and preference data in enhancing performance in complex tasks. Key insights can be categorized as follows: (1) *Data*: (a) *Scale & Quality*: From small-model resampling (e.g., OpenMathInstruct (Toshniwal et al., 2024)) to large-scale synthetic data (e.g., Qwen-2.5-MATH (Yang et al., 2024a)), datasets now include millions of samples. (b) *Efficiency*: Approaches like "less is more" alignment (e.g., LIMA (Zhou et al., 2023)) demonstrate that minimal, high-quality data can optimize pretrained capabilities. (2) *Optimization Framework*: (a) *Sampling Strategies*: Online RL techniques (e.g., DeepSeek V3 (DeepSeek-AI, 2024)) mitigate distributional shifts. (b) *Training Paradigms*: Multi-stage, collaborative optimization (e.g., Llama 3’s DPO iteration) improves model performance. (c) *Algorithms*: Advancements in PPO techniques, such as DPO and GRPO, focus on reducing parameter count and refining reward functions (e.g., PRIME (Cui et al., 2025)). These trends emphasize efficiency, generalization, and precision in unlocking LLMs’ reasoning potential.

Insight from LLM alignment The development of LLM alignment highlights three key insights and opportunities for improvement: (1) training efficiency—PPO-based methods require simultaneous loading of policy and reference models, slowing training; reference-free approaches like SimPO (Meng et al., 2024) could accelerate optimization by eliminating dependency on reference models, though their role in MLLM alignment needs deeper analysis. (2) overoptimization mitigation (Gao

et al., 2023; Rafailov et al., 2024)—DPO/RLHF risks reward hacking where proxy metrics improve while real-world performance degrades, exacerbated by biased or low-quality data. Solutions include diversifying training datasets, early stopping, and regularization to balance generalization. Addressing these challenges requires rethinking optimization architectures, robust data curation, and synergistic integration of RL paradigms.

MLLM as agents Combine the advanced reasoning capabilities of LLMs with multimodal perception—encompassing text, images, and audio—enabling cross-modal knowledge synthesis and task decomposition for complex real-world applications (Xi et al., 2023; Wang et al., 2024b; Ma et al., 2024b; Durante et al., 2024; Ma et al., 2024a). These capabilities position MLLMs as promising agents for various domains, such as autonomous driving and industrial robotics (Li et al., 2023b; Liu et al., 2024b). However, designing MLLMs as effective agents presents several unresolved challenges: (1) *Multi-agent Collaboration*: Lack of mature frameworks for multimodal communication (Ossowski et al., 2025), shared memory, and coordination in MLLM-based multi-agent systems. (2) *Robustness*: Vulnerability to adversarial attacks (e.g., image perturbations hijacking agent behavior (Wu et al., 2025)) in open environments, necessitating systematic robustness testing and defense mechanisms. (3) *Security*: Expanded attack surfaces across multimodal perception, reasoning, and memory modules, requiring comprehensive safeguards against privacy breaches and malicious hijacking (Yang et al., 2024b).

E Leveraging Visual Information for Alignment

For clarity, we use the following notation to represent the composition of current alignment data: preference data $\mathcal{D} = (x, \mathcal{I}, y_w, y_l)$, where x is the question, \mathcal{I} is the image, and y_w and y_l represent the winning and losing responses, respectively. In current research, three main approaches are employed to leverage visual information in order to enhance alignment performance, though each has its limitations:

1. **Using corrupted or irrelevant images as alignment phase negative samples.** Researchers create new images \mathcal{I}_{neg} and use $(y_w|x, \mathcal{I}_{neg})$ as a negative sample. This approach improves MLLM robustness to differ-

ent images and reduces hallucinations. However, visual negatives often rely on diffusion algorithms or image modifications that lack robust quality metrics, incurring high computational costs.

2. **Generating new questions and answers**

based on corrupted images. In this method, researchers create a new image \mathcal{I}_{neg} , use it to generate additional response y_{neg} , and then treat $(y_{neg}|x, \mathcal{I})$ as a negative sample. This method also essentially compares textual outputs, but it adds more variety to the textual comparison. However, the process of generating additional negative samples incurs extra computational overhead.

3. **Using cosine similarity metrics from models like CLIP to assess text-image matching.**

This approach uses a similarity score between the text and the image to filter data or as part of the reinforcement learning reward function. While this can help reduce data noise, the quality of the score depends on the evaluation model's quality, which may be subject to model bias.

Each of these methods plays a role in enhancing MLLM alignment with visual data, but they come with trade-offs in terms of efficiency, cost, and the potential for biases.

Method	Loss
Fact-RLHF	$\mathcal{L}_{\text{RLHF}} = -\mathbf{E}_{(\mathcal{I},x) \in \mathcal{D}, y \sim \pi_{\phi}(y \mathcal{I},x)} [r_{\theta}(\mathcal{I}, x, y) - \beta \cdot \mathbb{D}_{KL}(\pi_{\phi}(y \mathcal{I}, x) \parallel \pi^{\text{INIT}}(y \mathcal{I}, x))]$
SILKIE SIMA CLIP-DPO RLAIF-V 3D-CT-GPT++ MAVIS EMMOE xGen-MM(BLIP-3) LLaVA-NeXT-Interleave LLaVA-CRITIC SQuBa PPLLaVA HDPO SymDPO INTERACTIVECOT SAE-V TPR SAFEVIS CPCF LLaVA-Reasoner-DPO OmniAlign-V AIGI-Holmos IPA MM-IFEngine SPR C-DPO Modified DPO	$\mathcal{L}_{\text{dpo}} = -\mathbf{E}_{(\mathcal{I},x,y_w,y_l) \sim \mathcal{D}} [\log \sigma(\beta \log \frac{\pi_{\theta}(y_w \mathcal{I},x)}{\pi_{\text{ref}}(y_w \mathcal{I},x)}) - \beta \log \frac{\pi_{\theta}(y_l \mathcal{I},x)}{\pi_{\text{ref}}(y_l \mathcal{I},x)})]$
RLHF-V	$\mathcal{L}_{\text{Dense-dpo}} = -\mathbf{E}_{(\mathcal{I},x,y_w,y_l) \sim \mathcal{D}} [\mathbb{I}_{y_l \neq y_w} [\log \sigma(\beta \log \frac{\pi_{\theta}(y_w \mathcal{I},x)}{\pi_{\text{ref}}(y_w \mathcal{I},x)}) - \beta \log \frac{\pi_{\theta}(y_l \mathcal{I},x)}{\pi_{\text{ref}}(y_l \mathcal{I},x)})] + \mathbb{I}_{y_l \in y_w} [\gamma \log \sigma(\beta \log \frac{\pi_{\theta}(y_w \mathcal{I},x)}{\pi_{\text{ref}}(y_w \mathcal{I},x)}) - \beta \log \frac{\pi_{\theta}(y_l \mathcal{I},x)}{\pi_{\text{ref}}(y_l \mathcal{I},x)})]]$
F-DPO	$\mathcal{L}_{\text{Fine grained-dpo}} = -\mathbf{E}_{(\mathcal{I},x,y_w,y_l) \sim \mathcal{D}} [\log \sigma(\beta \log \frac{\pi_{\theta}(y_w \mathcal{I},x)}{\pi_{\text{ref}}(y_w \mathcal{I},x)}) - \log \sigma(\beta \log \frac{\pi_{\theta}(y_l \mathcal{I},x)}{\pi_{\text{ref}}(y_l \mathcal{I},x)})]$
HA-DPO & MSR-ViR	$\mathcal{L} = \mathcal{L}_{\text{dpo}} + \mathbf{E}_{(\mathcal{I},x,y) \sim \mathcal{D}_{\text{SFT}}} [-\log P(y \mathcal{I}, x; \pi_{\theta})]$
MIA-DPO	Loss : $\mathcal{L} = \mathcal{L}_{\text{dpo}} + \gamma \cdot \mathbf{E}_{(\mathcal{I},x,y_w,y_l) \sim \mathcal{D}} [-\log(y_w \mathcal{I}, x)]$
ChiP	$\mathcal{L} = \mathcal{L}_{\text{dpo}} + \mathcal{L}_{\text{Image dpo}} + \lambda \cdot \mathcal{L}_{\text{Sentence dpo}} + \gamma \cdot \mathbf{E}_{(\mathcal{I},x,y_w^{\text{Token}},y_l^{\text{Token}}) \sim \mathcal{D}_{\text{Token}}} [\beta \mathbb{D}_{\text{SeqKL}}[\pi_{\text{ref}}(y_w \mathcal{I}, x) \parallel \pi_{\theta}(y_w \mathcal{I}, x)] - \beta \mathbb{D}_{\text{SeqKL}}[\pi_{\text{ref}}(y_l \mathcal{I}, x) \parallel \pi_{\theta}(y_l \mathcal{I}, x)]]$
Image DPO	$\mathcal{L}_{\text{Image dpo}} = -\mathbf{E}_{(\mathcal{I}_w,\mathcal{I}_l,x,y_w)} [\log \sigma(\beta \log \frac{\pi_{\theta}(y_w \mathcal{I}_w,x)}{\pi_{\text{ref}}(y_w \mathcal{I}_w,x)}) - \beta \log \frac{\pi_{\theta}(y_l \mathcal{I}_l,x)}{\pi_{\text{ref}}(y_l \mathcal{I}_l,x)})]$
AdPO	$\mathcal{L} = -\mathbf{E}_{(\mathcal{I}_w,\mathcal{I}_l,x,y_w,y_l)} [\log \sigma(\beta \log \frac{\pi_{\theta}(y_w \mathcal{I}_w,x)}{\pi_{\text{ref}}(y_w \mathcal{I}_w,x)}) - \beta \log \frac{\pi_{\theta}(y_l \mathcal{I}_l,x)}{\pi_{\text{ref}}(y_l \mathcal{I}_l,x)})] + \sum_{t=1}^T \log \pi_{\theta}(y_w^t \mathcal{I}_l, x_t^{t-1})$
PHANTOM	$\mathcal{L} = \mathcal{L}_{\text{SFT}} - \mathbf{E}_{(\mathcal{I}_w,\mathcal{I}_l,x,y_w)} [\log \sigma(\frac{\beta}{ y_w } \log \pi_{\theta}(y_w \mathcal{I}_w, x)) - (\frac{\beta}{ y_l } \log \pi_{\theta}(y_l \mathcal{I}_l, x))]$
video-SALMONN 2	$\mathcal{L} = \mathcal{L}_{\text{dpo}} + \lambda \mathbf{E}_{(\mathcal{I},x,y_{\text{gt}}) \sim \mathcal{D}_{\text{gt}}} [\log \pi_{\theta}(y_{\text{gt}} \mathcal{I}, x)]$
Preference Optimization	$\mathcal{L} = \mathcal{L}_{\text{dpo}} + \lambda \mathbf{E}_{(\mathcal{I},x,y) \sim \mathcal{D}_{\text{reg}}} [\log \frac{\pi_{\theta}(y x)}{\pi_{\text{ref}}(y x)}]$
DAMA	$\mathcal{L} = -\mathbf{E}_{(\mathcal{I},x,y_w,y_l) \sim \mathcal{D}} [\log \sigma(\alpha \cdot \beta \log \frac{\pi_{\theta}(y_w \mathcal{I},x)}{\pi_{\text{ref}}(y_w \mathcal{I},x)}) - \alpha \cdot \beta \log \frac{\pi_{\theta}(y_l \mathcal{I},x)}{\pi_{\text{ref}}(y_l \mathcal{I},x)})]$
mDPO	$\mathcal{L} = \mathcal{L}_{\text{dpo}} + \mathbf{E}_{(\mathcal{I}_w,\mathcal{I}_l,x,y_w,y_l) \sim \mathcal{D}} [-\log \sigma(\beta \log \frac{\pi_{\theta}(y_w \mathcal{I}_w,x)}{\pi_{\text{ref}}(y_w \mathcal{I}_w,x)}) - \beta \log \frac{\pi_{\theta}(y_l \mathcal{I}_l,x)}{\pi_{\text{ref}}(y_l \mathcal{I}_l,x)})]$
LPOI	$-\log \sigma(\beta \log \frac{\pi_{\theta}(y_w \mathcal{I}_w,x)}{\pi_{\text{ref}}(y_w \mathcal{I}_w,x)}) - \delta$
MPO	$\mathcal{L} = \alpha_1 \cdot \mathcal{L}_{\text{dpo}} - \alpha_2 \cdot \mathbf{E}_{(\mathcal{I},x,y_w,y_l) \sim \mathcal{D}} [\log \sigma(\beta \log \frac{\pi_{\theta}(y_w \mathcal{I},x)}{\pi_{\text{ref}}(y_w \mathcal{I},x)}) - \delta] - \alpha_2 \cdot [\log \sigma(\beta \log \frac{\pi_{\theta}(y_l \mathcal{I},x)}{\pi_{\text{ref}}(y_l \mathcal{I},x)}) - \delta] - \alpha_3 \cdot [\log \pi_{\text{ref}}(y_w \mathcal{I},x)]$
CcDPO MMedPO MFPO HSCR IMG BPO	$\mathcal{L}_{\text{Image-Text DPO}} = \alpha_1 \cdot \mathcal{L}_{\text{DPO}} + \alpha_2 \cdot \mathcal{L}_{\text{Image DPO}}$
PanoDPO	$\mathcal{L}_{\text{Image-2-Text DPO}} = \mathcal{L}_{\text{Image-Text DPO}} + \alpha_3 \cdot -\mathbf{E}_{(\mathcal{I},x,y_w,y_l) \sim \mathcal{D}} [\log \sigma(\beta \log \frac{\pi_{\theta}(y_w \mathcal{I},x+c)}{\pi_{\text{ref}}(y_w \mathcal{I},x+c)}) - \beta \log \frac{\pi_{\theta}(y_l \mathcal{I},x)}{\pi_{\text{ref}}(y_l \mathcal{I},x)})]$
D ² PO PAR	$\alpha_1 \mathcal{L}_{\text{Action DPO}} + \alpha_2 \mathcal{L}_{\text{State DPO}}$ $\mathcal{L}_{\text{simPO}} = -\mathbf{E}_{(\mathcal{I},x,y_w,y_l) \sim \mathcal{D}} [\log \sigma(\beta \log \pi_{\theta}(y_w \mathcal{I}, x)) - \beta \log \pi_{\theta}(y_l \mathcal{I}, x)]$
DCD	$\mathcal{L}_{\text{DCD}} = -\mathbf{E}_{(\mathcal{I},x,y_w,y_l) \sim \mathcal{D}} [\log \pi_{\theta}(y_l x, \mathcal{I}_l) + \log \pi_{\theta}(y_w x, \mathcal{I}_w)]$
fDPO	$\mathcal{L}_{\text{fDPO}} = \beta_1 \cdot \mathcal{L}_{\text{DescDPO}} + \beta_2 \cdot \mathcal{L}_{\text{ReasonDPO}}$
Pair-DPO	$\mathcal{L}_{\text{Pair-dpo}} = -\mathbf{E}_{(\mathcal{I},x,y_w,y_l) \sim \mathcal{D}} [\log \sigma(\Delta \text{Und} \Delta \text{Gen})]$ $\Delta \text{Und} = \beta \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - \beta \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)}$, $\Delta \text{Gen} = \beta \log \frac{\pi_{\theta}(y_w x)}{\pi_{\text{ref}}(y_w x)} - \beta \log \frac{\pi_{\theta}(y_l x)}{\pi_{\text{ref}}(y_l x)}$

Table 1: Various preference optimization objectives given preference data $\mathcal{D} = (x, \mathcal{I}, y_w, y_l)$, where x is the question, \mathcal{I} is the Image, and y_w and y_l are winning and losing responses.

Dataset	Size	Categories	Response Model	Data Sources	Annotation Model
LLaVA-RLHF	10K	Hallucination	LLaVA-SFT	LLaVA-Instruct	Human
RLHF-V	1.4K	Hallucination	Muffin	UniMM-Chat	Human
VLFeedback	80K	Hallucination	12 Models	9 Datasets	GPT-4
CLIP-DPO	750K	Hallucination	MobileVLM-v2	12 Datasets	CLIP
M-HalDetect	16K	Hallucination	InstructBLIP	MS COCO	Human
HA-DPO	6K	Hallucination	3 Models	Visual Genome	GPT-4
SIMA	17K	Hallucination	LLaVA-1.5	LLaVA-Instruct	LLaVA-1.5
RLAIF-V	83K	Hallucination	3 Models	7 Datasets	2 Models
xGen-MM (BLIP-3)	62.6K	Hallucination	xGen-MM-4B	open-source	-
MIA-DPO	52K	Multi-Image	LLaVa-v1.5 & InternLM-XC 2.5	Not mentioned	Not mentioned
MAVIS	88K	Math	MAVIS-7B	Self-Constructed	GPT-4
EMMOE-100	10K	Embodied AI	Video-LLaVA	Self-Constructed	GPT-4
Image-DPO	60K	visual reasoning	Cambrian-8B & LLaVA-1.5	3 Datasets	Stable Diffusion
LLAVA-CRITIC	40.1K	Multiple tasks	LLaVA-OneVision	3 Datasets	LLaVA-OneVision
MMPR	3.25M	Reasoning	InternVL2-8B	Not mentioned	automate pipeline
video-SALMONN-o1	100K	Reasoning	Gemini-1.5-pro	Self-Constructed	GPT4
MMedPO	19 K	Medical	Med-LVLM & GPT-4o	Self-Constructed	GPT-4o
CPCF	50K	Image Referring	Fine-tuned Forret	Open Images	MLLM Itself
SHAREGPT-4O-REASONING	193K	Reasoning	Fine-tuned LLaVA-Next	9 Datasets	Golden Answer
OmniAlign-V	205K	Multiple tasks	LLaVA-Next & GPT-4o	Self-Constructed	automate pipeline
D ² PO	15K	Embodiment AI	automate pipeline	Self-Constructed	automate pipeline
Holmes-DPOSet	49K	Multiple tasks	Fine-tuned MLLMs	Self-Constructed	automate pipeline
Pick-a-Pic	851K	Generation	Self-Generated	Generative Models	
IPA	89K	Instruction Following	Qwen2VL-7B	15 Datasets	LLM-as-judge
MM-IFInstruct	23k	Instruction Following	Intern-VL-2.5 7B	automate pipeline	
PAR	5K	General	LLaVA-1.5 7B & Qwen-VL-chat-9B	LLaVA-80K	KL divergence
SPR	10K	Spatial Understanding	Ferret	Object365	automate pipeline
SENTINEL	8.6K	Hallucination	LLaVA-1.5	Visual Genome	automate pipeline
TAlt & PAlt	18K	alt-text generation	Gemini	Social Networks	automate pipeline
Chart2Code	23K	code generation from chart	GPT-4o	ReachQA	Reward Models
MMSafe-PO	5.6K	Safety	human & MLLM	Anthropic-HH	Human
Fin-APT	470	Financial Advisory Videos	Gemma2-9B	Youtube	GPT-4o
Pair-DPO	5K	Multimodal Understanding and Generation	Show-o	Journey-DB	automate pipeline
SafeVid-350K	350K	Video Understanding Safety	LLaVA-Next & GPT-4	Self-Constructed	automate pipeline
TPR	20K	hallucination	LLaVA-1.5-7B	7 datasets	LLaVA-Next-34B
MultiScope	42k	Multi-Image	3 datasets	LLaVA-OneVision	automated-pipeline

Table 2: Preference optimization dataset construction, including dataset, data size, categories: usage of the data, data sources, response model: the model to generate responses y_w and y_l by given image \mathcal{I} and prompt x , and annotation model: the model to annotate y_w and y_l .