

COGGUIDE: HUMAN-LIKE GUIDANCE FOR ZERO-SHOT OMNI-MODAL REASONING

Zhou-Peng Shou^{1,2}, Zhi-Qiang You¹, Fang Wang¹, Hai-Bo Liu^{3*}

¹NoDesk AI, Hangzhou, China

²Zhejiang University, Hangzhou, China

³Independent Researcher, Hangzhou, China

ABSTRACT

Addressing shortcut reliance and limited contextual understanding in cross-modal reasoning, we introduce a zero-shot omni-modal reasoning component inspired by human-like cognition, realized as a plug-and-play intent-sketch pipeline with three serial modules: Intent Perceiver, Policy Generator and Strategy Selector, that models an understand-plan-select cognitive process. By producing and filtering lightweight intent sketch strategies to guide reasoning, the method requires no parameter fine-tuning and enables cross-model transfer through in-context engineering. An information-theoretic analysis demonstrates that this process reduces conditional entropy and improves information utilization efficiency. Extensive experiments on IntentBench, WorldSense, and Daily-Omni confirm the method’s generality and robustness: the full three-module design consistently outperforms strong baselines across diverse reasoning engines and pipeline settings, with maximum gains of +9.51 pp and a relative improvement of 20.04%. These results demonstrate the practical value and portability of the proposed intent sketch reasoning component for zero-shot omni-modal reasoning.

Index Terms—LLM, Omni-Modal, Intent Sketch, Information Entropy

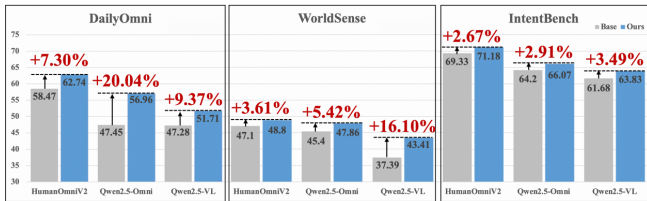


Fig 1: Relative Percentage Improvements over Baselines by the Proposed Method on IntentBench, WorldSense, and Daily-Omni

1. INTRODUCTION

With GPT-4 introducing visual inputs in 2023, LLMs began parsing images and performing visual reasoning [1], spurring rapid progress in building and applying multimodal models [2], [3]. Tasks like video question answering demand stronger dynamic-scene comprehension and temporal reasoning [4], motivating unified alignment and fusion of vision, language, and audio to enhance cross-modal reasoning [5]. As a result, models are expanding from images to video and audio toward “audiovisual omnipotent” architectures that enrich human-computer interaction [6].

However, even with massive parameters and multimodal pre-training, they continue to face challenges such as insufficient global contextual understanding and “shortcut” reasoning in complex tasks [7]. Models often over-rely on local or single-modality cues while

overlooking critical cross-modal information, leading to outputs that deviate from human intent [8]. Even with Chain-of-Thought prompts, multimodal LLMs remain weak at multi-step cross-modal reasoning [9]. Constructing reasoning chains via reinforcement learning may cause models to acquire “shortcut” strategies, thereby reducing generalization [10]. These phenomena indicate that relying solely on the model’s inherent reasoning ability and simple prompts is still insufficient to align the model with human intent.

To address this, recent work injects “intent” as a mediator between user queries and cross-modal evidence: some introduce explicit intent labels or scene purposes to constrain answers and reasoning scope [7]; others use instruction tuning and templated prompts to declare desired behavior and implicitly align goals [11]; still others, “intent-conditioned” retrieval-reasoning pipelines or agents drive evidence selection and reasoning steps by current intent [12]; and text-guided fusion aids multimodal intent understanding [13]. Increasingly, temporal and event structures in audio-video are exploited to infer latent intent and filter irrelevant cues [14], [15]. However, these approaches often depend on dense annotation and task-specific training, hindering zero-shot transfer [7], [16], or treat intent as static labels or prompt fragments without externalizing it into generable, assessable, and selectable strategies—limiting suppression of shortcuts and local biases [8], [11].

Motivated by recent “sketches of thought” approaches: the Sketch-of-Thought framework [17] maintains reasoning accuracy while reducing verbose intermediate reasoning, Machine Mental Imagery [18] injects latent visual “imagination” to aid reasoning in complex scenes. We propose a zero-shot omni-modal reasoning component inspired by human-like cognition. The component decomposes context through a plug-and-play intent-sketch pipeline: Intent Perceiver, Policy Generator and Strategy Selector, which infer textual intent from video and audio, generate candidate policies, and select an optimal strategy; a reasoning LLM then conditions on the selected strategy to produce the final answer via prompt composition. On challenging benchmarks—IntentBench [7], WorldSense [14] and Daily-Omni [15], using HumanOmniV2, Qwen2.5-Omni, Qwen2.5-VL, respectively, serving as a unified comparison. Our method surpasses the base model on all benchmarks (highest relative increase: 20.04%). Meanwhile, in terms of “effectiveness vs. cost” compared with training-based methods, we maintain zero training overhead and low latency, resulting in better overall cost and deployment timelines.

The main contributions: (1) a zero-shot omni-modal reasoning component that contains a plug-and-play intent-sketch pipeline, deliver immediate accuracy gains and cross-model generality; (2) an information-entropy analysis that explains how strategy prompts reduce decision uncertainty; (3) present comprehensive omni-modal experiments, including ablations and reasoning-engine swaps, that validate the effectiveness and practicality of in-context strategy prompting on the three benchmarks above.

*Corresponding author: Hai-Bo Liu (email: haiboliu2025@gmail.com)

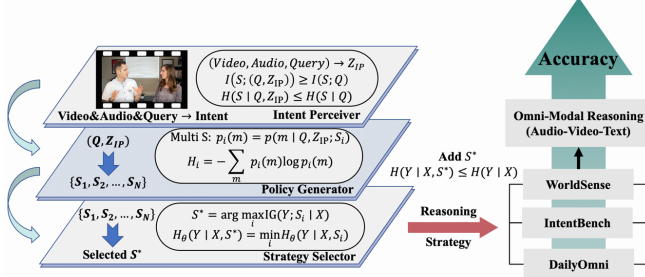


Fig 2: zero-shot omni-modal reasoning component: plug-and-play intent-sketch pipeline and information-entropy analysis

2. METHOD

The plug-and-play intent-sketch pipeline contains three serial modules: Intent Perceiver, Policy Generator and Strategy Selector, operating on omni-modal input $X = (V, A, Q)$ where V is video, A is audio, and Q is the textual query. The Intent Perceiver maps X to Z_{IP} to reduce generating uncertainty; the Policy Generator yields N semantically diverse policies; the Strategy Selector applies extra conditioning C , select S^* to drive the final reasoning. This design mirrors the human routine of understand-plan-select, yielding a clearer reasoning trajectory and improved answer generation.

2.1. Intent Perceiver: Intent Representation for Strategy Generation

Given $X = (V, A, Q)$, the module outputs an intent representation $Z_{IP} = f(X)$ that fuses cross-modal cues with the query. From an information-theoretic view, using Z_{IP} to condition strategy generation provides nonnegative information gain:

$$I(S; (Q, Z_{IP})) - I(S; Q) = I(S; Z_{IP} | Q) \geq 0 \quad (1)$$

equivalently $H(S | Q, Z_{IP}) \leq H(S | Q)$. By focusing on answer-relevant omni-modal evidence, Z_{IP} reduces ambiguity at problem understanding and effectively lowers the initial conditional uncertainty $H(Y | X)$.

2.2. Policy Generator: Omni-Policy Generation Based on Semantic Entropy

Conditioned on (Q, Z_{IP}) , a frozen LLM (prompted as a policy provider) generates N candidate sketches $\{S_1, \dots, S_N\}$ that specify reasoning lines without producing the answer. Let \mathcal{M} denote semantic equivalence classes over policies and $p(m | Q, Z_{IP})$ the posterior over classes. We define the semantic entropy

$$H_{sem}(S | Q, Z_{IP}) = - \sum_{m \in \mathcal{M}} p(m | Q, Z_{IP}) \log p(m | Q, Z_{IP}) \quad (2)$$

and the gain brought by intent conditioning:

$$I(S; Z_{IP} | Q) = H_{sem}(S | Q) - H_{sem}(S | Q, Z_{IP}) \geq 0 \quad (3)$$

To balance single-sketch clarity and set-level coverage, for each candidate S_i let $p_i(m) = p(m | Q, Z_{IP}; S_i)$ and

$$H_i = - \sum_m p_i(m) \log p_i(m) \quad (4)$$

For the mixture $p(m) = \frac{1}{N} \sum_{i=1}^N p_i(m)$ with entropy $H(\bar{p}) = - \sum_m p(m) \log p(m)$, we optimize

$$\max_{S_1, \dots, S_N} H(\bar{p}) - \alpha \frac{1}{N} \sum_{i=1}^N H_{sem}(S_i | Q, Z_{IP}) + \gamma \text{Div}(S_1, \dots, S_N) \quad (5)$$

$\alpha, \gamma > 0$ are weights, and $\text{Div}(\cdot)$ can be defined via pairwise distances in semantic similarity between candidates to encourage complementary policies with different emphases, such as “evidence-first,” “temporal/causal-first,” or “cross-modal-alignment-first”.

2.3. Strategy Selector: Strategy Selection Based on Minimum Conditional Entropy/Bayesian Risk

Given $\{S_1, \dots, S_N\}$, the model (prompted as a strategy evaluator) selects S^* based on the X and task-specific conditioning C to align the choice with real-world constraints. Let $p_\theta(Y | X, S_i, C)$ be the answer posterior under S_i with extra conditioning C . In implementation, we obtain $p_\theta(Y | X, S_i, C)$ by prompting the LLM with (X, S_i, C) and computing answer-slot likelihoods, so the added C induces strategy-dependent posteriors. Under 0-1 loss:

$$S^* = \arg \min_i R_{0-1}(S_i) = \arg \max_i \max_Y p_\theta(Y | X, S_i, C) \quad (6)$$

Selection maximizes information gain:

$$IG(Y; S_i | X) = H_\theta(Y | X) - H_\theta(Y | X, S_i, C) \quad (7)$$

$$S^* = \arg \max_i IG(Y; S_i | X) = \arg \min_i H_\theta(Y | X, S_i, C) \quad (8)$$

Because candidates are complementary:

$$H_\theta(Y | X, S^*) = \min_i H_\theta(Y | X, S_i) \leq \frac{1}{N} \sum_i H_\theta(Y | X, S_i, C) \quad (9)$$

In practice, exploring a diverse and complementary set of strategies increases the chance that at least one S_i substantially reduces uncertainty, making lower overall uncertainty more likely. This implies a lower expected conditional entropy—via a Fano-type bound—a lower achievable error bound; in particular, reducing $H_\theta(Y | X, S^*, C)$ lowers an upper bound on the minimum error rate under 0-1 loss. Consequently, subsequent reasoning based on S^* integrates omni-modal evidence along a path of higher confidence, enhancing the reliability of answer generation.

2.4. A Unified Information-Theoretic Framework of Uncertainty Reduction

This section formalizes our pipeline as an information-theoretic mechanism for uncertainty reduction. We begin with the basic theorem of conditioning:

$$H(X | Y) \leq H(X) \quad (10)$$

Treating the strategy S^* as an observable planning variable, C denotes the extra task conditioning used to choose S^* , we model the answer posterior as $p(Y | X, I, S, C)$. Formally, by the conditioning reduces entropy theorem, we obtain the following monotone contraction chain:

$$H(Y | X) \geq H(Y | X, I) \geq H(Y | X, I, S) \geq H(Y | X, I, S, C) \quad (11)$$

This captures how intent I , the multi-policy set S , and the selected strategy S^* progressively reduce the conditional uncertainty of Y .

From the perspective of the Data Processing Inequality (DPI), if $X \rightarrow I \rightarrow S \rightarrow S^*$ forms a Markov chain, then we have:

$$I(X; S^*) < I(X; S) < I(X; I) < I(X; X) \quad (12)$$

This implies that any representation derived from X must lose some information about X itself. However, in inference, we do not replace X with S^* ; instead, we use S^* as auxiliary information conditioned on X :

$$H(Y | X, S^*) < H(Y | X) \quad (13)$$

with strict inequality when $I(Y; S^* | X) > 0$. This reduction occurs because S^* , although derived from X , is constructed through a process that incorporates task-specific guidance (e.g., evidence-first, causal-first, cross-modal alignment). These human-informed design choices embed additional relevant cues into S^* , enabling it to provide new information about Y that is not directly accessible from X alone. Thus, conditioning on S^* enhances prediction accuracy, leading to the entropy reduction in (13) and contributing to the overall uncertainty contraction in (11).

Algorithm 1 Pseudocode for intent sketch reasoning component

Input $X = (V, A, Q)$, k ; **Output** \hat{y}, R
 $I \leftarrow \text{IntentPerceiver}(X)$
 $S \leftarrow \text{PolicyGenerator}(I, Q, k) = \{s_1, \dots, s_k\}$
function POSTERIOR_AND_ENTROPY(X, s)
 $\hat{P}_Y \leftarrow \text{ReasoningEngine.posterior}(Y | X, s)$
 $\hat{H} \leftarrow \text{Entropy}(\hat{P}_Y)$
return (\hat{P}_Y, \hat{H})
 $\text{best_H} \leftarrow +\infty$; $\text{best_s} \leftarrow \emptyset$; $\text{post_best} \leftarrow \emptyset$
for each s **in** S **do**
 $\hat{P}_Y, \hat{H} \leftarrow \text{POSTERIOR_AND_ENTROPY}(X, s)$
if $\hat{H} < \text{best_H}$:
 $\text{best_H} \leftarrow \hat{H}$; $\text{best_s} \leftarrow s$; $\text{post_best} \leftarrow \hat{P}_Y$
end for
 $s^* \leftarrow \text{best_s}$
 $(R, \hat{y}) \leftarrow \text{ReasoningEngine.solve_with_strategy}(X, s^*)$
return \hat{y}, R

3. EXPERIMENTS AND RESULTS

3.1. Experimental Settings

We evaluate the method on three omni-modal reasoning benchmarks: IntentBench (omni-intent understanding), WorldSense (audio-video collaborative analysis), and Daily-Omni (daily-life scenarios). Accuracy (%) is the metric. We compare baseline models (no strategy prompts) with our approach that adds three plug-in modules—Intent Perceiver, Policy Generator, and Strategy Selector—whose output is passed to the reasoning engine (Table I).

All experiments follow a zero-shot setting (no fine-tuning; zero-shot prompts only). We use three reasoning engines: HumanOmniV2, Qwen2.5-Omni, and Qwen2.5-VL. The pipeline modules are instantiated with four pretrained LLMs: closed-source GPT-4o and Doubao-Seed-1.6, and open-source GLM-4.5 and Qwen3. To quantify module contributions, we conduct ablations (Table II). Each setting is evaluated over the Cartesian product of the four pipeline LMs and the three reasoning engines, with all other hyperparameters and inference configurations held constant for fair comparison. Through these combinations we measure the overall gain of the integrated system and the contribution of each module.

TABLE I

Summary of Models, Roles, and Scales (“a/b” denotes total / activated parameters, for Mixture-of-Experts models)

Model	Role	Parameter Scale
-------	------	-----------------

HumanOmniV2[7]	Reasoning Engine	7B
Qwen2.5-Omni	Reasoning Engine	7B
Qwen2.5-VL	Reasoning Engine	7B
GPT-4o	Policy Generator /Strategy Selector	Large Closed-Source Model
GLM-4.5	Policy Generator /Strategy Selector	355B/32B
Doubao-Seed-1.6	Policy Generator /Strategy Selector	Large Closed-Source Model
Qwen3	Policy Generator /Strategy Selector	235B/22B
Qwen2.5-VL-32B	Intent Perceiver	32B
GLM-4.5V	Intent Perceiver	106B/12B

TABLE II

Experimental Configurations: Three-Module Settings

Exp ID	Intent Perceiver	Policy Generation	Strategy Selection
CG_Qwen	Qwen2.5-VL-32B	3	On
CG_GLM	GLM-4.5V	3	On
Abl_NI	No	3	On
Abl_SP	Qwen2.5-VL-32B	1	On
BaseLine	No	No	No

Note: CG_Qwen: Full, use Qwen as Intent model; CG_GLM: Full, use GLM as Intent model; Abl_NI: Remove the Intent module; Abl_SP: Change policy generation to single policy; BaseLine: No front-end pipeline.

3.2. Main Results

Table III report detailed results on three benchmarks, respectively. The best results are sometimes achieved by CG_Qwen and sometimes by the CG_GLM, but both surpass “no intent” (Abl_NI) and “single policy” (Abl_SP). In addition, both Abl_NI and Abl_SP outperform the baseline, confirming the effectiveness of each module. the maximum gain corresponds to +9.51 pp (a 20.04% relative improvement). Importantly, the three-module scheme consistently outperforms the corresponding baselines across all combinations of the four pipelines and three reasoning engines, regardless of whether the pipeline LLMs are closed-source or open-source, this indicates that our method offers strong portability and plug-and-play characteristics.

TABLE III

Accuracy (%) of different reasoning models combined with Pipeline module models on IntentBench (IN), WorldSense (WO) and Daily-Omni (DA) datasets.

Pipeline	Reasoning Model (Baseline)	CG_Qwen			CG_GLM			Abl_NI			Abl_SP		
		IN	WO	DA	IN	WO	DA	IN	WO	DA	IN	WO	DA
GPT-4o	Human OmniV2	70.86	48.8	60.23	70.47	48.55	59.9	70.45	47.79	58.56	70.09	48.14	59.31
GLM-4.5	(IN:69.33)	71.07	48.17	62.74	70.87	48.7	61.24	70.51	48.01	59.31	70.27	47.89	60.74
Doubao-Seed-1.6	(WO:47.1)	70.92	48.23	62.66	70.72	48.2	61.07	70.06	48.14	59.9	69.74	47.64	60.9
Qwen3	(DA:58.47) [7]	71.18	48.36	62.49	70.9	48.36	61.32	70.82	48.3	61.24	69.96	47.92	60.15

GPT-4o	Qwen2.5-Omni	65.95	47.13	55.56	66.07	47.67	55.64	65.82	47.01	55.47	64.99	46.75	50.38
GLM-4.5	(IN: 64.2)	65.51	47.57	54.05	65.86	47.38	54.22	65.31	47.23	53.38	65.31	46.31	50.88
Doubao-Seed-1.6	(WO: 45.4)	65.67	47.45	54.89	65.6	47.04	52.38	65.45	46.94	51.88	65.3	46.69	50.54
Qwen3	(DA: 47.45)	65.67	47.86	56.9	65.83	47.79	56.96	65.45	47.7	56.81	65.46	46.22	51.88
GPT-4o		62.72	43.1	49.71	62.75	43.1	49.96	62.64	42.97	49.62	62.14	41.87	49.12
GLM-4.5	Qwen2.5-VL	63.12	43.41	51.55	63.25	42.88	51.71	62.69	42.4	51.38	62.4	41.93	50.96
Doubao-Seed-1.6	(IN: 61.68)	63.2	42.21	50.79	63.02	42.12	50.63	62.9	41.93	50.46	62.09	41.33	50.35
Qwen3	(WO: 37.39)	63.81	43.06	51.63	63.83	43.25	51.55	63.78	42.91	51.46	63.39	41.83	50.54
	(DA: 47.28)												

Note: The first column lists pipeline models. Each pipeline is paired with four reasoning models, shown in the subsequent rows. Baseline values (shown in bold within parentheses) indicate each reasoning model’s standalone performance on the corresponding dataset.

3.3. Ablation Study and Analyses

We ablate the three modules: Intent Perceiver (IP), Policy Generator (PG), and Strategy Selector (SS), with results in Table III. Removing any single component degrades accuracy. For clarity, we denote Abl-NI as removing IP (i.e., no Z_{IP}), and Abl-SP as using a single policy ($N = 1$), which disables PG’s diversity and makes SS degenerate to selecting the only sketch. For Daily-Omni (HumanOmniV2 \times GLM-4.5), the full system reaches 62.74%; Abl-NI drops to 59.31% (−3.43 pp), and Abl-SP to 60.74% (−2.00 pp). Similar patterns hold elsewhere: on WorldSense (Qwen2.5-VL \times GLM-4.5) full 43.41%, no-IP 42.40% (−1.01 pp), single-policy 41.93% (−1.48 pp); on IntentBench (HumanOmniV2 \times Qwen3) full 71.18%, no-IP 70.82% (−0.36 pp), single-policy 69.96% (−1.22 pp).

Interpretation. Without the intent provided by the Intent Perceiver, the model may miss crucial multi-modal cues, typically resulting in a modest performance drop; with Abl-SP ($N = 1$), the strategy set collapses and coverage over semantic classes M shrinks, which aligns with the increase of semantic uncertainty in Eq. (2)-(3). Note that Abl-SP does not remove SS; rather, the selector trivially chooses the only available sketch. Therefore, the observed decline under Abl-SP should be attributed to the loss of PG-induced diversity (and the $\text{Div}(\cdot)$ term in Eq. (5)). Among the ablations, the impact of policy generation is larger: the single-policy (Abl-SP) degradation tends to exceed the no-IP case, indicating that multi-path thinking is intrinsically valuable. This mirrors the human “understand–plan–select” pattern, where careful multi-path simulation precedes a final decision. The selector’s role is captured by minimum conditional entropy/Bayesian risk (Eq. (6)-(9)); when $N > 1$ and SS actively scores candidate sketches, uncertainty is further reduced in line with the contraction in Eq. (11)-(13). From the pipeline perspective, Qwen3-based pipelines exhibit the highest win rate across settings; given that most reasoning engines here are from the Qwen family, this “same-lineage” pairing likely improves adherence to strategy guidance and reduces mismatch between policy sketches and the executor. From the reasoning-model perspective, Omni-model outperform VL-model, suggesting that access to audio cues materially improves disambiguation and intent grounding.

Our approach is model and framework agnostic. Across reasoning engines (HumanOmniV2, Qwen2.5-Omni, Qwen2.5-VL) and pipelines (GPT-4o / GLM-4.5 / Doubao-Seed-1.6 / Qwen3), improvements are consistent. Gains are larger for base models lacking post-training (e.g., WorldSense: Qwen2.5-VL 37.39% \rightarrow 43.41%, +6.02 pp, CG Qwen, Pipeline=GLM-4.5; HumanOmniV2 47.10% \rightarrow 48.80%, +1.70 pp, Pipeline=GPT-4o). On Daily-Omni,

HumanOmniV2 improves 58.47% \rightarrow 62.74% (+4.27 pp) with GLM-4.5; replacing GLM-4.5 by the smaller Qwen3 still yields +4.02 pp. These results indicate a plug-and-play component that reliably enhances reasoning across open/closed models and scales, demonstrating cross-model, cross-platform, and cross-scenario robustness.

4. CONCLUSION

This paper presents a zero-shot omni-modal reasoning component that implements a plug-and-play, intent-sketch pipeline: Intent Perceiver, Policy Generator, and Strategy Selector, which uses context injection to enhance reasoning performance without fine-tuning and generalizes across modalities, reasoning engines, and platforms. From an information-theoretic view, Policy Generator produces complementary candidates conditioned on intent, Strategy Selector selects among candidate strategies via prompt-based evaluation under additional task-specific conditioning, effectively contracting the answer posterior and reducing decision uncertainty, explaining the stable gains and improved interpretability without training overhead.

On IntentBench, WorldSense, and Daily-Omni, the full three-module scheme consistently outperforms baselines across all Pipeline/engine combinations, with maximum gains of +9.51 pp (relative improvement of 20.04%). Ablations show complementary roles: PG+SS is the primary source of improvement, while explicit intent is especially beneficial for video/audio-dependent tasks. The method delivers stable cross-platform gains for diverse reasoning engines and Pipeline models, underscoring its plug-and-play nature and strong transferability. Overall, the intent sketch offers a lightweight, effective paradigm for improving the alignment, robustness, and interpretability of complex cross-modal reasoning.

5. REFERENCES

- [1] S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg, H. Nori, H. Palangi, M. T. Ribeiro, and Y. Zhang, "Sparks of Artificial General Intelligence: Early experiments with GPT-4," *arXiv preprint arXiv:2303.12712*, 2023.
- [2] D. Driess, F. Xia, M. S. M. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu, W. Huang, Y. Chebotar, P. Sermanet, D. Duckworth, S. Levine, V. Vanhoucke, K. Hausman, M. Toussaint, K. Greff, A. Zeng, I. Mordatch, and P. Florence, "PaLM-E: An Embodied Multimodal Language Model," in *ICML*, 2023, pp. 8469-8488.
- [3] Z. Yan, Z. Li, Y. He, C. Wang, K. Li, X. Li, X. Zeng, Z. Wang, Y. Wang, Y. Qiao, L. Wang, and Y. Wang, "Task Preference Optimization: Improving Multimodal Large Language Models with Vision Task Alignment," in *CVPR*, 2025, pp. 29880-29892.
- [4] J. Li, P. Wei, W. Han, and L. Fan, "IntentQA: Context-aware Video Intent Reasoning," in *ICCV*, 2023, pp. 11963-11974.
- [5] M. Ma, Z. Yu, Y. Ma, G. Li and Z. Yang, "EventLens: Enhancing Visual Commonsense Reasoning by Leveraging Event-Aware Pretraining and Cross-modal Linking," in *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2025, pp. 1-5.
- [6] D. Zhang, Y. Yu, J. Dong, C. Li, D. Su, C. Chu, and D. Yu, "MM-LLMs: Recent Advances in MultiModal Large Language Models," in *ACL*, 2024, pp. 12401-12430.
- [7] Q. Yang, S. Yao, W. Chen, S. Fu, D. Bai, J. Zhao, B. Sun, B. Yin, X. Wei, and J. Zhou, "HumanOmniV2: From Understanding to Omni-Modal Reasoning with Context," *arXiv preprint arXiv:2506.21277*, 2025.
- [8] X. Zheng, C. Liao, Y. Fu, K. Lei, Y. Lyu, L. Jiang, B. Ren, J. Chen, J. Wang, C. Li, L. Zhang, D. P. Paudel, X. Huang, Y.-G. Jiang, N. Sebe, D. Tao, L. V. Gool, and X. Hu, "MLLMs are Deeply Affected by Modality Bias," *arXiv preprint arXiv:2505.18657*, 2025.
- [9] Y. Hao, J. Gu, H. W. Wang, L. Li, Z. Yang, L. Wang, and Y. Cheng, "Can MLLMs Reason in Multimodality? EMMA: An Enhanced MultiModal Reasoning Benchmark," in *ICML*, 2025.
- [10] J. Xia, Y. Zang, P. Gao, Y. Li, and K. Zhou, "Visionary-R1: Mitigating Shortcuts in Visual Reasoning with Reinforcement Learning," *arXiv preprint arXiv:2505.14677*, 2025.
- [11] Z. Hu, A. Iscen, C. Sun, Z. Wang, K.-W. Chang, Y. Sun, C. Schmid, D. A. Ross, and A. Fathi, "REVEAL: Retrieval-Augmented Visual-Language Pre-Training with Multi-Source Multimodal Knowledge Memory," in *CVPR*, 2023, pp. 23369-23379.
- [12] S. Ravi, A. Chinchure, L. Sigal, R. Liao, and V. Shwartz, "VLC-BERT: Visual Question Answering with Contextualized Commonsense Knowledge," in *WACV*, 2023, pp. 1155-1165.
- [13] Y. Zhang, B. Chen, H. Ye, Z. Gao, T. Wan and L. Lan, "Text-guided Multimodal Fusion for the Multimodal Emotion and Intent Joint Understanding," in *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2025, pp. 1-2.
- [14] J. Hong, S. Yan, J. Cai, X. Jiang, Y. Hu, and W. Xie, "WorldSense: Evaluating Real-world Omnimodal Understanding for Multimodal LLMs," *arXiv preprint arXiv:2502.04326*, 2025.
- [15] Z. Zhou, R. Wang, and Z. Wu, "Daily-Omni: Towards Audio-Visual Reasoning with Temporal Alignment across Modalities," *arXiv preprint arXiv:2505.17862*, 2025.
- [16] Z. Zhang, A. Zhang, M. Li, H. Zhao, G. Karypis, and A. Smola, "Multimodal Chain-of-Thought Reasoning in Language Models," in *TMLR*, 2024.
- [17] S. A. Aytes, J. Baek, and S. J. Hwang, "Sketch-of-Thought: Efficient LLM Reasoning with Adaptive Cognitive-Inspired Sketching," *arXiv preprint arXiv:2503.05179*, 2025.
- [18] Z. Yang, X. Yu, D. Chen, M. Shen, and C. Gan, "Machine Mental Imagery: Empower Multimodal Reasoning with Latent Visual Tokens," *arXiv preprint arXiv:2506.17218*, 2025.