# VISUAL PUNS FROM IDIOMS: AN ITERATIVE LLM-T2IM-MLLM FRAMEWORK

*Kelaiti Xiao[1,2]   Liang Yang[1]   Dongyu Zhang[1]   PAERHATI Tulajiang[1,2]   Hongfei Lin[1*]*

[1]Dalian University of Technology, Dalian, China
[2]Xinjiang Normal University, Urumqi, China

## ABSTRACT

We study idiom-based visual puns—images that align an idiom's literal and figurative meanings—and present an iterative framework that coordinates a large language model (LLM), a text-to-image model (T2IM), and a multimodal LLM (MLLM) for automatic generation and evaluation. Given an idiom, the system iteratively (i) generates detailed visual prompts, (ii) synthesizes an image, (iii) infers the idiom from the image, and (iv) refines the prompt until recognition succeeds or a step limit is reached. Using 1,000 idioms as inputs, we synthesize a corresponding dataset of visual pun images with paired prompts, enabling benchmarking of both generation and understanding. Experiments across 10 LLMs, 10 MLLMs, and one T2IM (Qwen-Image) show that MLLM choice is the primary performance driver: GPT achieves the highest accuracies, Gemini follows, and the best open-source MLLM (Gemma) is competitive with some closed models. On the LLM side, Claude attains the strongest average performance for prompt generation. Code is available at `https://github.com/xkt88/VisualPun`.

***Index Terms—*** Visual Puns, Idiom Understanding, Multimodal Large Language Models, Benchmark dataset

## 1. INTRODUCTION

Visual puns are a sophisticated form of multimodal communication that combines linguistic wordplay with visual representation, requiring simultaneous processing of literal and metaphorical meanings [1, 2]. Creating and interpreting such puns demands resolving semantic ambiguities, aligning visual and linguistic cues, and exercising creative reasoning [3]. These properties make visual puns valuable benchmarks for assessing AI models' creative and interpretive capabilities [4].

Given recent advances in T2IMs [5–7], a natural direction is to leverage them to automatically synthesize visual-pun datasets at scale. State-of-the-art T2IMs can faithfully render specific scenes when provided with sufficiently detailed prompts [8–10]. However, they lack the deeper linguistic reasoning needed to capture figurative intent, so metaphorical prompts are often misinterpreted [11]. An effective strategy is to leverage an LLM to decompose an idiom's literal and metaphorical facets and translate them into concrete visual directives for the T2IM, as illustrated in Fig. 1. Even with such LLM-guided prompting, the community still lacks standardized resources and reliable procedures.
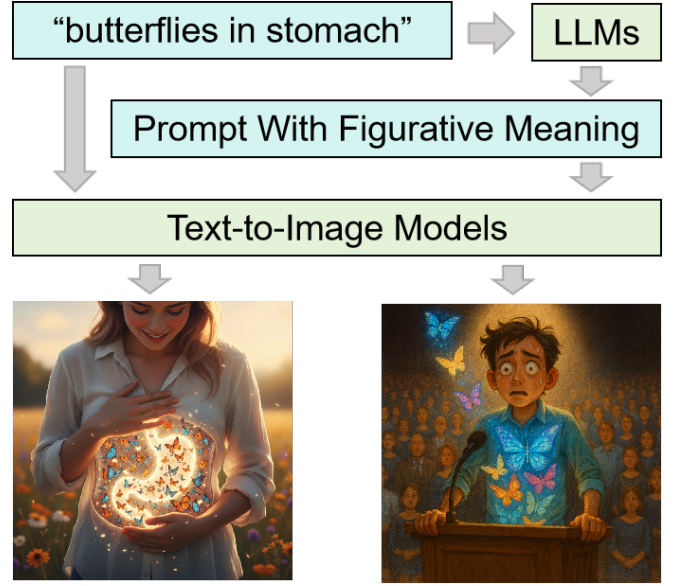


**Fig. 1**. An LLM bridges literal and metaphorical meanings of an idiom, generating prompts that guide a T2IM to create visual puns capturing both nervous feelings and literal butterflies.

Despite rapid progress in text-to-image generation and multimodal modeling [12–15], two gaps remain: (i) there is no large-scale public benchmark specifically targeting idiom-based visual puns with synthesized images [16]; existing resources focus on visual metaphors or rebus art and differ in scope [17–19]; and (ii) one-shot generations are unreliable due to hallucinations and semantic ambiguity [20], motivating iterative self-improvement loops that detect, diagnose, and correct mismatches between intent and output. To address these challenges, we contribute:

- An iterative LLM–T2IM–MLLM pipeline that decomposes idioms into literal and figurative cues, generates/refines prompts, synthesizes images, and stops when the idiom is recognized or a step cap is reached.
- Using 1,000 English idioms, we produce a one-to-one dataset of visual-pun images with paired prompts, enabling benchmarking of generation and understanding.
- A large-scale evaluation across 10 LLMs and 10 MLLMs showing MLLM choice dominates performance, while T2IM choice matters less with detailed prompts.
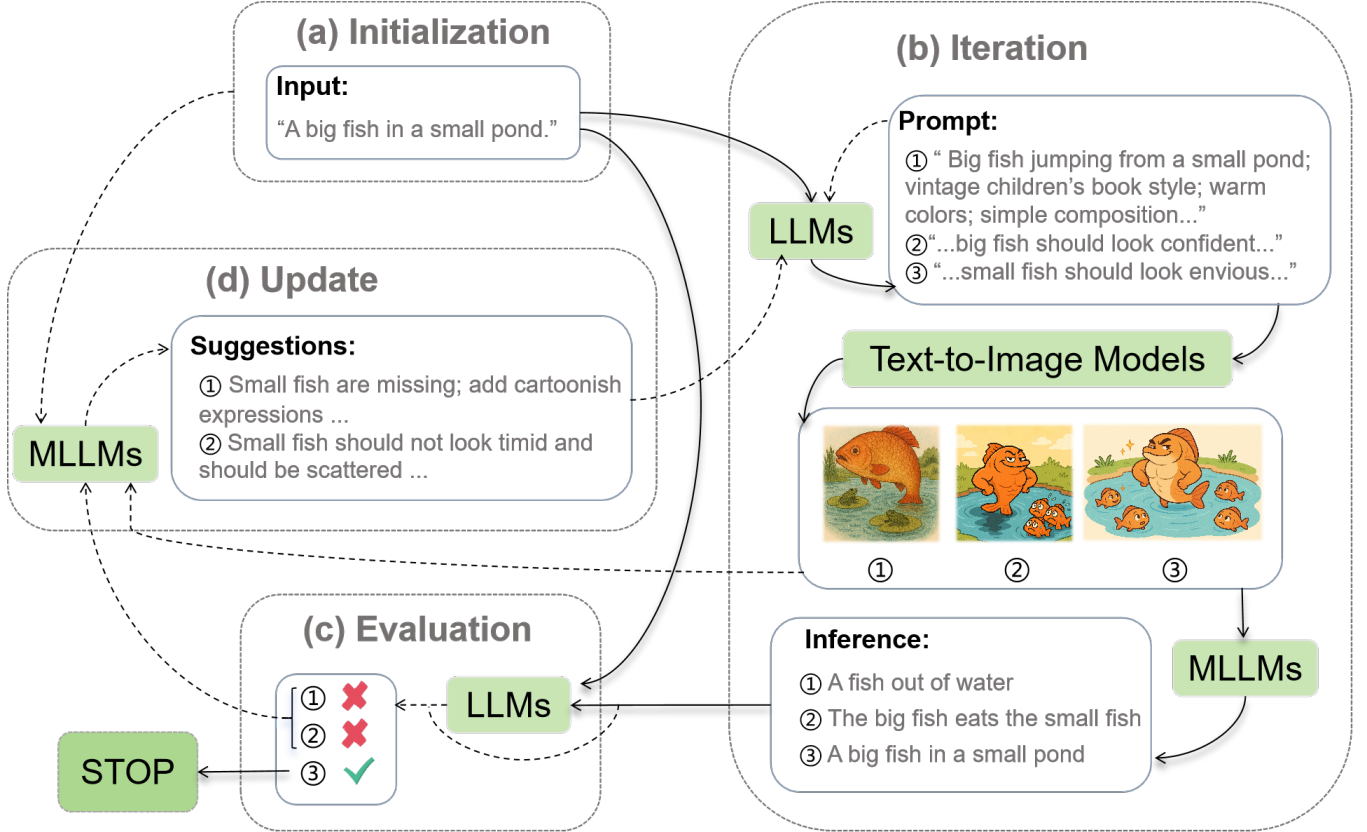
**Fig. 2**. Iterative pipeline for visual pun generation: the LLM crafts prompts, the T2IM synthesizes images, the MLLM infers the idiom, and the LLM evaluates and updates; solid arrows denote the current iteration, dashed arrows reference the previous iteration

## 2. METHODOLOGY

Our pipeline (Fig. 2) iterates four modules to render visual puns that convey literal and figurative meanings. Given an idiom $I_{\text{input}}$, an LLM proposes a visual prompt $P_t$ conditioned on prior prompt $P_{t-1}$ and edit suggestions $U_{t-1}$. A T2IM synthesizes an image $G_t$; an MLLM infers the idiom $R_t$ from $G_t$. An LLM judge checks semantic equivalence and issues a control signal; on mismatch, the MLLM supplies edits $U_t$ to refine $P_t$. We stop on a match or after $t=5$ iterations.

### 2.1. Initialization

Given an idiom corpus $\mathcal{I} = \{I_1, \ldots, I_n\}$, select one target idiom:

$$I_{\text{input}} = I_j, \quad j \in \{1, \ldots, n\}. \tag{1}$$

Initialize $P_0 = \varepsilon$ (empty prompt) and $U_0 = \emptyset$ (no suggestions).

### 2.2. Iteration

At iteration $t$:

$$P_t = \text{LLM}_{\text{prompt}}(I_{\text{input}}, U_{t-1}, P_{t-1}), \tag{2}$$
$$G_t = \text{T2IM}(P_t), \tag{3}$$
$$R_t = \text{MLLM}_{\text{infer}}(G_t), \tag{4}$$

where $P_t$ is the textual prompt, $G_t$ the synthesized image, and $R_t$ the inferred idiom (top-1 string).

### 2.3. Evaluation

An LLM judge tests semantic equivalence (with canonicalization of idiom form):

$$M_t = \text{LLM}_{\text{eval}}(R_t, I_{\text{input}}) \in \{\text{true}, \text{false}\}. \tag{5}$$

The control signal is

$$C_t = \begin{cases} \text{STOP}, & M_t = \text{true or } t \geq 5, \\ \text{CONTINUE}, & \text{otherwise}. \end{cases} \tag{6}$$

### 2.4. Update

If $C_t = \text{CONTINUE}$, we produce targeted refinements for the next iteration:

$$U_t = \text{MLLM}_{\text{update}}(R_t, G_t, I_{\text{input}}), \tag{7}$$

where $U_t$ lists concrete edits to $P_t$ (e.g., missing objects, composition, emphasis).

## 3. EXPERIMENTS

### 3.1. Experimental Setup

Table 2 lists the evaluated models: 10 LLMs and 10 MLLMs, with 6 models used in both roles (prompt generation and visual understanding). Closed-source models were accessed via

**Table 1**. Idiom recognition accuracy (%) for the iterative pipeline across 10 MLLMs (rows) and 10 LLMs (columns) on 1,000 idioms; the bottom row reports column-wise averages.

| MLLMs | | Close-sourced LLMs | | | | | Open-sourced LLMs | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | GPT | Gemini | Claude | Grok | Doubao | DeepSeek | GPT-OSS | Llama | GLM-4.5 | Qwen3 |
| Close-Sourced | GPT | 76.9 | 73.7 | 79.8 | 69.5 | 67.3 | 70.1 | 71.1 | 64.8 | 65.9 | 68.7 |
| | Gemini | 71.8 | 69.5 | 74.8 | 65.1 | 63.1 | 65.7 | 66.7 | 60.8 | 61.6 | 64.4 |
| | Claude | 59.7 | 57.6 | 61.6 | 54.4 | 52.6 | 54.9 | 55.5 | 50.8 | 51.6 | 53.8 |
| | Grok | 58.8 | 56.5 | 61.8 | 52.9 | 51.6 | 53.6 | 54.5 | 49.6 | 50.2 | 52.5 |
| | Doubao | 55.6 | 53.4 | 58.2 | 50.4 | 49.5 | 51.0 | 51.6 | 47.1 | 48.0 | 50.0 |
| Open-Sourced | Llama | 45.9 | 44.1 | 47.8 | 41.6 | 40.3 | 42.1 | 42.6 | 38.8 | 39.5 | 41.2 |
| | GLM-4.5V | 48.7 | 46.8 | 50.7 | 44.2 | 42.8 | 44.5 | 45.2 | 41.2 | 41.9 | 43.7 |
| | Qwen2.5 | 51.3 | 49.3 | 53.4 | 46.5 | 45.1 | 46.9 | 47.6 | 43.4 | 44.1 | 46.0 |
| | Gemma | 56.1 | 54.0 | 58.1 | 50.8 | 49.2 | 51.3 | 52.0 | 47.4 | 47.9 | 50.4 |
| | Mistral | 28.2 | 27.1 | 29.4 | 25.6 | 24.8 | 25.9 | 26.2 | 24.1 | 24.3 | 25.3 |
| Average | | 55.3 | 53.2 | 57.6 | 50.1 | 48.6 | 50.6 | 51.3 | 46.8 | 47.5 | 49.6 |

**Table 2**. Overview of evaluated models by provider and release date; bold marks models used in both LLM and MLLM tasks.

| Models | Company | Released Date |
|---|---|---|
| *Close Sourced (M)LLMs* | | |
| **GPT-5** [21] | OpenAI | August 7, 2025 |
| **Gemini-2.5-flash** [22] | Google | June 17, 2025 |
| **Claude-Sonnet-4** [23] | Anthropic | May 22, 2025 |
| **Grok-4** [24] | xAI | July 9, 2025 |
| **Doubao-Seed-1.6** [25] | ByteDance | June 11, 2025 |
| *Open Sourced (M)LLMs* | | |
| **Llama-4-Maverick** [26] | Meta | April 5, 2025 |
| DeepSeek-V3.1 [27] | DeepSeek | August 21, 2025 |
| GPT-OSS-120b [28] | OpenAI | August 5, 2025 |
| GLM-4.5 [29] | Z.ai | July 28, 2025 |
| Qwen3-32B [30] | Alibaba | April 29, 2025 |
| GLM-4.5V [31] | Z.ai | August 11, 2025 |
| Qwen2.5-VL-32B [32] | Alibaba | February 28, 2025 |
| Gemma-3-27b [33] | Google | March 12, 2025 |
| Mistral-Small-3.2 [34] | Mistralai | June 10, 2025 |
| *Text-to-Image Model* | | |
| Qwen-Image [35] | Alibaba | August 4, 2025 |

Poe API[1]; open-source models were accessed via the Deep-Infra API[2].

We use a single T2IM, Qwen-Image [35], with fixed settings, accessed through Poe API. Images were generated at a fixed resolution of $1024 \times 1024$. As evidenced in Sec. 3.4, when prompts are sufficiently detailed, different T2IMs exhibit comparable semantic fidelity. Therefore, we fix the T2IM to Qwen-Image to isolate LLM/MLLM effects and reduce variance and cost.

All experiments use 1,000 English idioms. The pipeline runs up to 5 iterations per idiom, generating one image per iteration, and stops early when the MLLM's top-1 inferred idiom matches the target after canonicalization. Accuracy is the proportion of idioms recognized under this criterion.

All prompts, hyperparameters, and other settings are unified across LLMs and MLLMs; for exact values and scripts, please refer to the released code, prompts, and images.

### 3.2. Main Results and Analysis

We report top-1 idiom recognition accuracy on 1,000 idioms (up to 5 iterations; fixed T2IM). Table 1 varies MLLMs by rows and LLMs by columns; the bottom row is the column-wise average across MLLMs.

**MLLM choice dominates performance.** GPT MLLM is best across all LLM partners (64.8–79.8%), followed by Gemini (60.8–74.8%). Claude and Grok form a mid-tier (mid-50s to low-60s), with Doubao slightly lower (high-40s to high-50s). Among open-source MLLMs, Gemma leads (47.4–58.1%), approaching closed-source Doubao on several columns; Qwen2.5 and GLM-4.5V cluster in the mid-40s to low-50s, Llama in the low-40s, and Mistral trails (24–29%). The best–worst MLLM gap at a fixed LLM can exceed 50 points (e.g., with Claude as LLM: 79.8 vs. 29.4), underscoring that visual understanding is the primary variance source; the MLLM ordering is largely stable across LLM partners.

**LLM effects are smaller but non-negligible.** Column-wise averages span 46.8–57.6%, with Claude strongest for prompt generation (57.6%), followed by GPT (55.3%) and Gemini (53.2%). While careful linguistic prompting helps, the narrower column range versus row spreads indicates the chief bottleneck lies in MLLM visual reasoning. The best observed pairing is GPT (MLLM) with Claude (LLM) at 79.8%.

**Open-source alternatives are competitive under cost constraints.** As MLLMs, Gemma (up to 58.1%) nears closed-source Doubao; as LLMs, GPT-OSS averages 51.3%, exceeding Doubao (48.6%) and narrowing the gap to top closed-source LLMs.

### 3.3. Ablation Study

We ablate three configurations on 1,000 idioms using the same T2IM (Qwen-Image, $1024 \times 1024$) and protocol as Sec. 3.1; accuracy is top-1 idiom match after canonicalization with early stopping (max 5 iterations).

**Table 3**. Ablation study showing incremental accuracy improvements from baseline T2IM generation through LLM-augmented prompting and iterative refinements (Claude as fixed LLM).

| Configuration | Close-sourced MLLMs | | | | |
|---|---|---|---|---|---|
| | GPT | Gemini | Claude | Grok | Doubao |
| T2IM | 52.3 | 45.2 | 37.5 | 36.8 | 32.3 |
| +LLM | 67.6 (+15.3) | 59.7 (+14.5) | 49.5 (+12.0) | 49.8 (+13.0) | 45.5 (+13.2) |
| Ours(1 update) | 76.2 (+8.6) | 68.8 (+9.1) | 57.5 (+8.0) | 58.3 (+8.5) | 54.3 (+8.8) |
| Ours(2 updates) | 79.3 (+3.1) | 74.5 (+5.7) | 61.0 (+3.5) | 61.5 (+3.2) | 57.8 (+3.5) |
| Ours(3 updates) | 79.7 (+0.4) | 74.8 (+0.3) | 61.5 (+0.5) | 61.8 (+0.3) | 58.2 (+0.4) |
| Ours(4 updates) | 79.8 (+0.1) | 74.8 (+0.0) | 61.6 (+0.1) | 61.8 (+0.0) | 58.2 (+0.0) |
| Configuration | Open-sourced MLLMs | | | | |
| | Llama | GLM-4.5V | Qwen2.5 | Gemma | Mistral |
| T2IM | 28.3 | 26.3 | 30.2 | 29.3 | 16.2 |
| +LLM | 38.3 (+10.0) | 38.8 (+12.5) | 42.3 (+12.1) | 44.1 (+14.8) | 23.5 (+7.3) |
| Ours(1 update) | 44.8 (+6.5) | 47.0 (+8.2) | 49.8 (+7.5) | 53.6 (+9.5) | 27.5 (+4.0) |
| Ours(2 updates) | 47.6 (+2.8) | 50.5 (+3.5) | 53.0 (+3.2) | 57.6 (+4.0) | 29.3 (+1.8) |
| Ours(3 updates) | 47.8 (+0.2) | 50.7 (+0.2) | 53.4 (+0.4) | 58.0 (+0.4) | 29.4 (+0.1) |
| Ours(4 updates) | 47.8 (+0.0) | 50.7 (+0.0) | 53.4 (+0.0) | 58.1 (+0.1) | 29.4 (+0.0) |



GPT-Image-1(0.1$)  DALL-E-3(0.03$)  Gemini-2.5-Image(0.02$)  Imagen-4(0.02$)  Qwen-Image(0.0132$)  Seedream-3(0.02$)
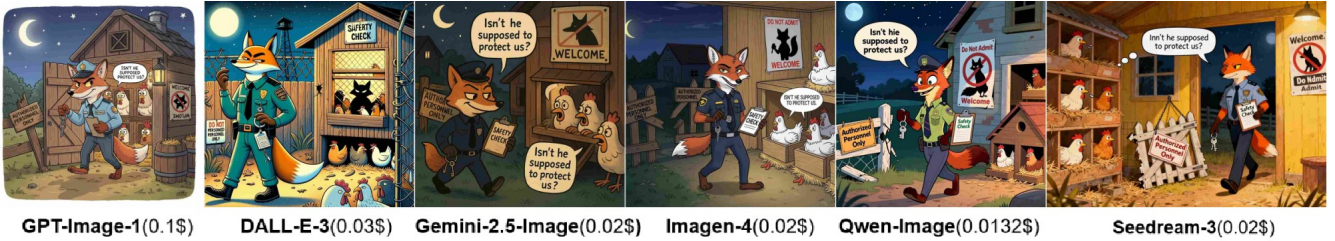
**Fig. 3**. Comparative visualization of the idiom "fox in a henhouse" generated by different T2IMs using identical prompting. All models were given the prompt:"*Cartoon night scene at a small farm:a sly,sharp-eyed fox in asecurity guard uniform holds a keyring and clipboard labeled"Safety Check,"confidently strolling into a cozy chicken coop.Nervous hens peek from nesting boxes,one whispering to another in a speech bubble,"Isn't he supposed to protect us?""*"

**Configurations.** (i) *T2IM-only*: the idiom string is used directly as the prompt; no LLM is involved, and the MLLMs performs recognition. (ii) *+LLM*: one-shot LLM-generated prompt (Claude fixed for prompting); no updates. (iii) *Ours (k updates)*: iterative prompt refinement with $k=1\ldots4$; Claude produces edits each round; the MLLMs handles recognition. All other settings are fixed; only the presence/number of updates varies.

**Results (Table 3).** T2IM-only yields 16.2–52.3% across MLLMs, indicating direct generation under-specifies figurative intent. Adding an LLM prompt improves accuracy by +7.3–15.3 points on every MLLM, showing the benefit of linguistic decomposition. Iterative refinement adds gains: the first update contributes +4.0–+9.5 points; later updates have diminishing returns, with negligible improvements by the 4th update, indicating convergence within 3–4 iterations.

**Takeaway.** With the T2IM held constant, performance gains are primarily attributable to (a) introducing linguistic guidance and (b) a small number of targeted updates, beyond which additional iterations bring little benefit.

### 3.4. Case Study

Using an identical LLM-crafted prompt for "fox in a henhouse" (Fig. 3; $1024 \times 1024$), several state-of-the-art T2IMs produced images that our MLLM consistently mapped to the target idiom despite stylistic differences. We observed the same pattern on 50 additional idioms across 6 T2IMs; for brevity, we release the selection protocol, idiom list, prompts, and images in the repository. Together with the ablations in Sec. 3.3—which attribute most gains to LLM guidance and a small number of iterative updates—this case study suggests that, under detailed prompting, T2IM choice is secondary. Accordingly, we fix Qwen-Image as the T2IM for all experiments.

### 4. CONCLUSION

We introduced an iterative framework for generating idiom-based visual puns and released a large-scale dataset of 1,000 idioms with paired prompts and images for benchmarking multimodal generation and understanding; resources are publicly available as cited in the abstract. Experiments varying 10 LLMs and 10 MLLMs with a fixed T2IM show that MLLM choice is the principal driver of recognition accuracy, while LLMs have smaller but consistent effects; most gains arise within 2–3 refinement rounds. Limitations include reliance on a single T2IM configuration and automatic, MLLM-based evaluation; future work will expand T2IM diversity and incorporate human studies and cross-lingual idioms. We hope these resources support more reliable assessment and progress in creative multimodal reasoning.

## 5. REFERENCES

[1] Jiwan Chung, Seungwon Lim, Jaehyun Jeon, et al., "Can visual language models resolve textual ambiguity with visual cues? let visual puns tell you!," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 2024, pp. 2452–2469.

[2] Yiran Rex Ma, Shan Huang, Yuting Xu, et al., "Pun2pun: Benchmarking llms on textual-visual chinese-english pun translation via pragmatics model and linguistic reasoning," in *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics*, 2025, pp. 331–354.

[3] Christian F Hempelmann and Andrea Samson, "Visual puns and verbal puns: Descriptive analogy or false analogy," *New Approaches to the Linguistics of Humour. Galati*, 2007.

[4] Yuan Yuan, Zhaojian Li, and Bin Zhao, "A survey of multimodal learning: Methods, applications, and future," *ACM Computing Surveys*, vol. 57, pp. 1 – 34, 2025.

[5] Zhijie Wang, Yuheng Huang, Da Song, et al., "Promptcharm: Text-to-image generation through multi-modal prompting and refinement," in *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, 2024, pp. 1–21.

[6] Lukas Höllein, Aljaž Božič, Norman Müller, et al., "Viewdiff: 3d-consistent image generation with text-to-image models," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024, pp. 5043–5052.

[7] Tong Liu, Zhixin Lai, Jiawen Wang, et al., "Multimodal pragmatic jailbreak on text-to-image models," in *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics*, July 2025, pp. 4681–4720.

[8] Chenshuang Zhang, Chaoning Zhang, Mengchun Zhang, et al., "Text-to-image diffusion models in generative ai: A survey," *arXiv preprint arXiv:2303.07909*, 2023.

[9] Pu Cao, Feng Zhou, Qing Song, and Lu Yang, "Controllable generation with text-to-image diffusion models: A survey," *arXiv preprint arXiv:2403.04279*, 2024.

[10] Rong-Cheng Tu, Zi-Ao Ma, Tian Lan, et al., "Automatic evaluation for text-to-image generation: Task-decomposed framework, distilled training, and meta-evaluation benchmark," in *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics*, 2025, pp. 22340–22361.

[11] Linhao Zhang, Jintao Liu, Li Jin, et al., "Gome: Grounding-based metaphor binding with conceptual elaboration for figurative language illustration," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 2024, pp. 18500–18510.

[12] Youwei Liang, Junfeng He, Gang Li, et al., "Rich human feedback for text-to-image generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 19401–19411.

[13] Junsong Chen, Chongjian Ge, Enze Xie, et al., "Pixart-$\sigma$: Weak-to-strong training of diffusion transformer for 4k text-to-image generation," in *European Conference on Computer Vision*. Springer, 2024, pp. 74–91.

[14] Fengxiang Bie, Yibo Yang, Zhongzhu Zhou, et al., "Renaissance: A survey into ai text-to-image generation in the era of large model," *IEEE transactions on pattern analysis and machine intelligence*, 2024.

[15] Hongji Yang, Yucheng Zhou, Wencheng Han, et al., "Self-rewarding large vision-language models for optimizing prompts in text-to-image generation," in *Findings of the Association for Computational Linguistics: ACL 2025*. July 2025, pp. 7332–7349, Association for Computational Linguistics.

[16] Jian Li, Weiheng Lu, Hao Fei, et al., "A survey on benchmarks of multimodal large language models," *arXiv preprint arXiv:2408.08632*, 2024.

[17] Tuo Zhang, Tiantian Feng, Yibin Ni, et al., "Creating a lens of chinese culture: A multimodal dataset for chinese pun rebus art understanding," *arXiv preprint arXiv:2406.10318*, 2024.

[18] Tuhin Chakrabarty, Arkadiy Saakyan, Olivia Winn, et al., "I spy a metaphor: Large language models and diffusion models co-create visual metaphors," in *Findings of the Association for Computational Linguistics: ACL 2023*, 2023, pp. 7370–7388.

[19] Chang Su, Xingyue Wang, Shupin Liu, et al., "Efficient visual metaphor image generation based on metaphor understanding," *Neural Processing Letters*, vol. 56, no. 3, pp. 150, 2024.

[20] Lei Huang, Weijiang Yu, Weitao Ma, et al., "A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions," *ACM Transactions on Information Systems*, vol. 43, no. 2, pp. 1–55, 2025.

[21] Shansong Wang, Mingzhe Hu, Qiang Li, et al., "Capabilities of gpt-5 on multimodal medical reasoning," *arXiv preprint arXiv:2508.08224*, 2025.

[22] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, et al., "Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities," *arXiv preprint arXiv:2507.06261*, 2025.

[23] Cindy Nguyen, Daniel Carrion, and Mohamed K Badawy, "Comparative performance of anthropic claude and openai gpt models in basic radiological imaging tasks," *Journal of Medical Imaging and Radiation Oncology*, 2025.

[24] xAI, "Grok 4," Online. Available: https://x.ai, Accessed: Sep. 16, 2025.

[25] ByteDance Seed Team, "Doubao-seed-1.6," Online. Available: https://seed.bytedance.com, Accessed: Sep. 16, 2025.

[26] Meta AI, "Llama 4 Maverick 17B 128E Original," Online. Available: https://huggingface.co/meta-llama/Llama-4-Maverick-17B-128E-Original, 2025, Accessed: Sep. 16, 2025.

[27] DeepSeek AI, "DeepSeek-V3.1," Online. Available: https://huggingface.co/deepseek-ai/DeepSeek-V3.1, 2025, Accessed: Sep. 18, 2025.

[28] Sandhini Agarwal, Lama Ahmad, Jason Ai, et al., "gpt-oss-120b & gpt-oss-20b model card," *arXiv preprint arXiv:2508.10925*, 2025.

[29] Aohan Zeng, Xin Lv, Qinkai Zheng, et al., "Glm-4.5: Agentic, reasoning, and coding (arc) foundation models," *arXiv preprint arXiv:2508.06471*, 2025.

[30] An Yang, Anfeng Li, Baosong Yang, et al., "Qwen3 technical report," *arXiv preprint arXiv:2505.09388*, 2025.

[31] Z.ai, "GLM-4.5V," Online. Available: https://huggingface.co/zai-org/GLM-4.5V, 2025, Accessed: Sep. 17, 2025.

[32] Shuai Bai, Keqin Chen, Xuejing Liu, et al., "Qwen2. 5-vl technical report," *arXiv preprint arXiv:2502.13923*, 2025.

[33] Gemma Team, Aishwarya Kamath, Johan Ferret, et al., "Gemma 3 technical report," *arXiv preprint arXiv:2503.19786*, 2025.

[34] Mistral AI, "Mistral Small 3.2 24B Instruct 2506," Online. Available: https://huggingface.co/mistralai/Mistral-Small-3.2-24B-Instruct-2506, 2025, Accessed: Sep. 16, 2025.

[35] Chenfei Wu, Jiahao Li, Jingren Zhou, et al., "Qwen-image technical report," *arXiv preprint arXiv:2508.02324*, 2025.