

VISTA: A Test-Time Self-Improving Video Generation Agent

Supplementary Material

6. Additional Discussions

VISTA’s Workflow. Figure 5 illustrates an overview of our proposed framework.

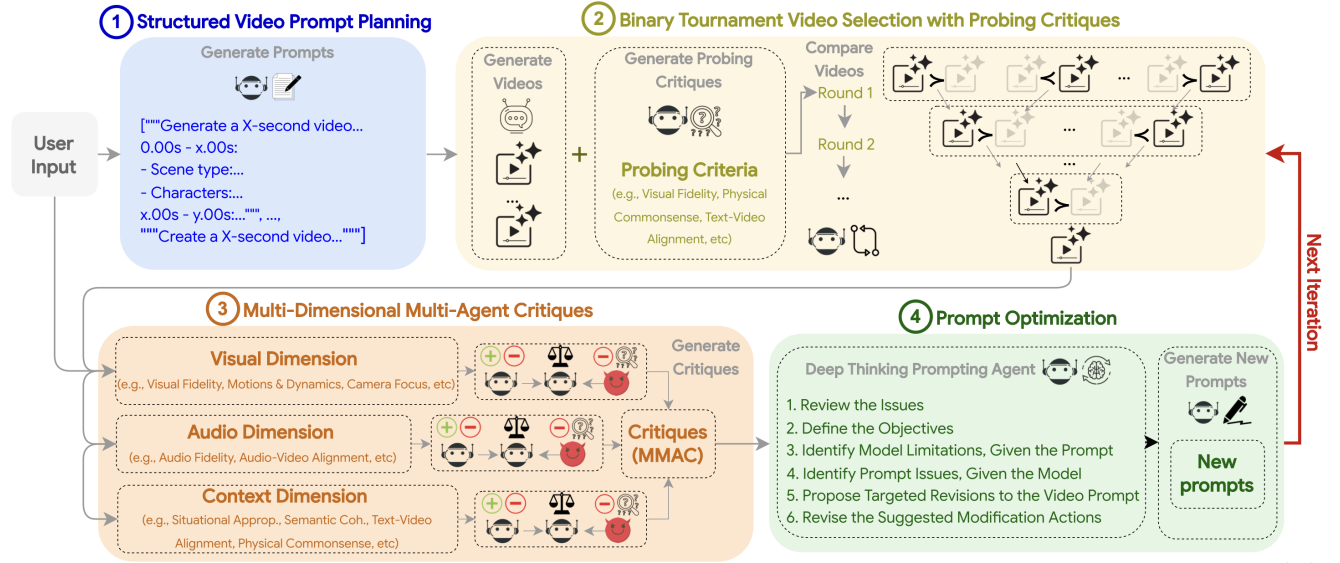


Figure 5. The workflow of our proposed multi-agent framework, VISTA. : MLLM Agent; : Adversarial MLLM Agent; : Video Generation Agent.

Cost Analysis. Figure 6 provides our cost-performance analyses.

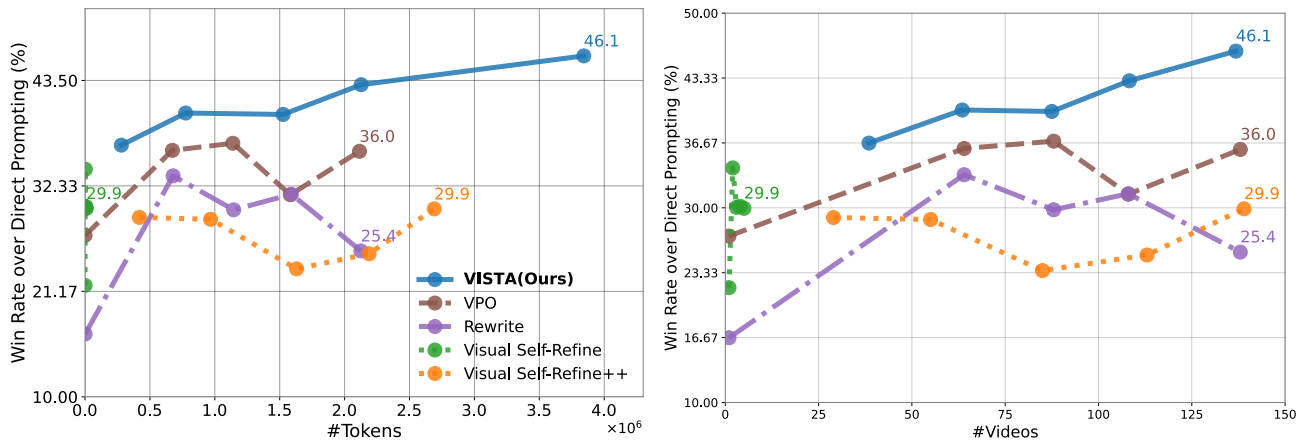


Figure 6. Cost analysis. Left: total token consumption, including both input and output tokens per iteration. Right: number of newly sampled videos per iteration. Results are averaged over two datasets. Tokens for video generation are unavailable and thus excluded.

Customizing VISTA. VISTA provides flexibility for users to define or adjust its behaviors. Both the selection metrics \mathcal{M}_{user}^S (Step 2), the critique metrics \mathcal{M}_{user}^C (Step 3), and all steps’ constraints can be customized. For example, the constraints in Step 1 can be removed to encourage model being more creative, leading to more creative videos generated from user prompt; \mathcal{M}_{user}^S can prioritize aspects that reflect subjective preferences, such as *color grading fidelity* or *emotional impact*,

emphasizing the users’ vision or affective goals. In addition, \mathcal{M}_{user}^C can focus on more fine-grained behaviors like the subtle fluidity of character gestures in the visual dimension. Further exploration of user-customized metrics and constraints, particularly those capturing unique aesthetic or narrative nuances, is an exciting direction to bring video generation closer to truly personalized creative video generation.

Latency Analysis. In practice, most of VISTA’s steps can be parallelized via batch inference. Let’s denote X as the latency of an MLLM call, and Y as of the video generator, where $X \ll Y$. Direct Prompting takes Y . VISTA’s Step 1 takes X , Step 2 takes X to get critiques while Y to generate videos in parallel. Tournament selection takes at most $X \times \lceil \log_2 n \rceil$ where n is #videos generated. MMAC takes $2 \times X$ (X for Normal and Negative judges, X for Meta judge), and prompt optimization takes $2 \times X$. So, one iteration takes $\approx Y + (\lceil \log_2 n \rceil + 6)X \approx Y$ with small n . Also, VISTA’s hyperparameters can be optimized: Fig. 4 shows that $n = 8$ already achieves good performance and we may not need $n \geq 10$ and 20 iterations. With optimal hardware, VISTA’s one iteration is expected to be comparable to DP.

Limitations. Despite the notable performance gains achieved by VISTA, several limitations remain, revealing promising avenues for future work. Firstly, our evaluation relies primarily on MLLMs and automated metrics, which may introduce systematic biases or fail to capture aspects of video quality that humans prioritize. While we include human evaluation and cross-validate with multiple automated approaches, the comprehensive human evaluation remains prohibitively expensive that affect the entire field. Secondly, while our critique mechanism provides a configurable interface, the current default metrics reflect certain assumptions about video quality that may not generalize across different cultural contexts, creative styles, or user preferences. Customizing VISTA’s metrics to better reflect user-specific or domain-specific preferences could enhance its adaptability and robustness. Lastly, VISTA requires both MLLMs and T2V models with strong instruction-following and reasoning capabilities to function effectively. As such models continue to improve, we expect this limitation to diminish.

7. Additional Results

7.1. Evaluations on Conventional Video and Audio Generation Metrics

Method	Subject Consistency	Background Consistency	Temporal Flickering	Motion Smoothness	Dynamic Degree	Aesthetic Quality	Imaging Quality	Temporal Style	Audio Noisiness	Audio Discontinuity	Audio Coloration	CLIP Score	Inception Score
DP	89.89	94.39	97.82	99.23	75.95	61.86	64.42	7.88	1.74	2.04	1.65	0.310	1.053
VSR	89.33	93.53	97.79	99.26	64.56	63.45	65.53	9.26	1.73	2.13	1.64	0.309	1.082
VSR++	87.96	93.53	97.88	99.12	74.68	60.68	63.06	9.25	1.65	2.03	1.56	0.310	1.078
Rewrite	89.09	93.79	97.59	99.17	77.22	62.52	62.58	8.57	1.64	1.99	1.53	0.310	1.085
VPO	86.74	92.66	97.76	99.15	77.22	61.17	64.01	8.03	1.70	1.97	1.59	0.311	1.039
VISTA	89.95	92.89	97.82	98.94	89.87	64.53	65.89	9.63	1.88	2.19	1.62	0.358	1.101

Table 4. **Single-scene:** Evaluation results using VBench’s any-video evaluation metrics for visual quality, NISQA metrics for audio quality, and CLIP-Score for text-video alignment.

Method	Subject Consistency	Background Consistency	Temporal Flickering	Motion Smoothness	Dynamic Degree	Aesthetic Quality	Imaging Quality	Temporal Style	Audio Noisiness	Audio Discontinuity	Audio Coloration	CLIP Score	Inception Score
DP	79.28	85.27	97.99	99.14	72.15	47.39	67.19	6.54	2.24	2.62	2.28	0.285	1.11
VSR	76.41	83.50	98.30	99.15	53.16	47.74	65.86	9.53	2.24	2.64	2.27	0.286	1.09
VSR++	78.52	85.43	98.33	99.11	60.75	47.96	65.68	6.80	2.14	2.53	2.13	0.288	1.08
Rewrite	77.85	86.08	98.02	99.08	67.09	50.56	65.62	8.73	1.93	2.38	2.05	0.290	1.11
VPO	77.08	83.80	98.15	99.13	58.23	49.08	67.04	8.32	2.25	2.67	2.26	0.288	1.09
VISTA	79.45	83.95	98.20	99.16	75.18	50.00	68.87	10.09	2.30	2.69	2.21	0.299	1.15

Table 5. **Multi-scene:** Evaluation results using VBench’s any-video evaluation metrics for visual quality, NISQA metrics for audio quality, and CLIP-Score for text-video alignment.

7.2. Human Evaluation Results Details

Table 6 presents our human evaluation results across annotators.

Ann.	VISTA	VSR(++)	Ann.	Single-scene VISTA	VSR++	Multi-scene VISTA	VSR	Ann.	Visual	Audio
Ann. 4	3.06	2.76	Ann. 1	76%	24%	72%	28%	Ann. 1	3.90 / 3.34	4.08 / 3.22
Ann. 5	4.44	4.02	Ann. 2	68%	32%	68%	32%	Ann. 2	3.72 / 3.50	3.64 / 3.14
Ann. 6	3.84	3.22	Ann. 3	64%	36%	60%	40%	Ann. 3	3.69 / 3.33	3.68 / 3.28
			Ann. 4	62%	38%	66%	34%			
			Ann. 5	66%	34%	62%	38%			

(a) Self-improvement scores (over 5).

(b) Win rate of VISTA vs. best baselines.

(c) Visual and audio scores (over 5).

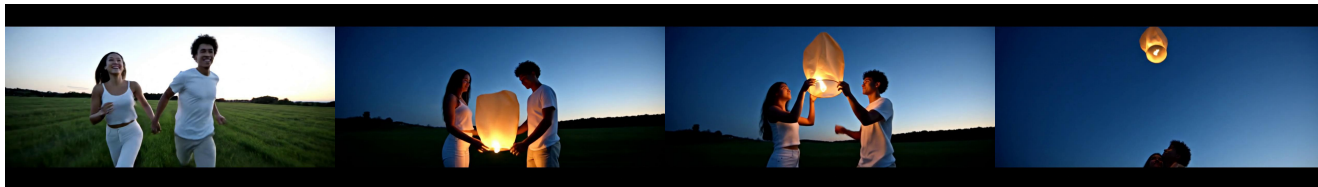
Table 6. Human evaluation results.



Prompt: A spaceship entering hyperdrive, stars streaking past as it accelerates.

Self-Refine: Overall, the generated video is highly successful in fulfilling the prompt...

VISTA (Motions and Dynamics): While the Normal Judge praises the smoothness of the ship’s motion and dynamic star streaking, the Negative Judge correctly identifies a major directional flaw: **the spaceship moves vertically, which conflicts with viewer expectations of horizontal acceleration**. Additionally, the Negative Judge points out **the lack of micro-dynamics (e.g., rotational drift, buildup phases) and unrealistic exhaust behavior, which diminish the believability of motion**...



Prompt: The couple runs hand in hand to release a sky lantern, then watches it drift upward into the night sky, carried by the wind with the stars shining above.

Self-Refine: The video successfully portrays the actions of the couple running and releasing a lantern...

VISTA (Contextual Suitability): The scene struggles significantly with internal logical consistency... **an abrupt and jarring shift from a bright, late-afternoon setting to a deep, artificial blue night sky** without any visual cues for time passing... **(Visual Characters):** While the two characters... their portrayal **lacks emotional depth and genuine interaction**. Their expressions, particularly during the lantern release, appear somewhat static and posed, missing an opportunity to convey a deeper sense of wonder...

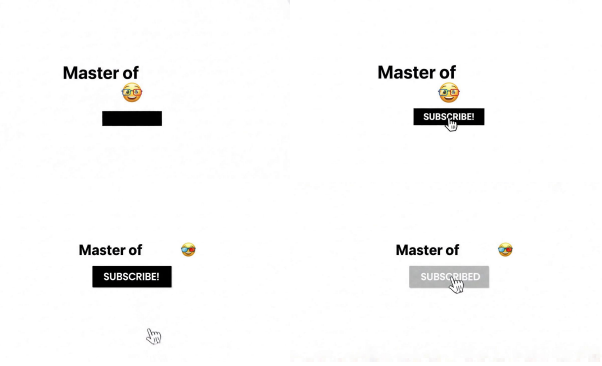
Table 7. VISTA delivers high-level critiques on visual fidelity, contextual consistency, and emotional expression, surfacing nuanced flaws that conventional LLM critiques often miss but are noticeable to human judges.

7.3. Why Does VISTA Work? Case Studies

VISTA Can Provide Human-Like Critiques (Step 3). VISTA provides intelligent, multi-faceted critiques that is often overlooked by conventional MLLM critiques, yet is intuitively recognized by humans. As shown in **Table 7**, VISTA can identify subtle visual flaws, such as a spaceship moving vertically instead of horizontally and movement against a static background. It also offers high-level reasoning context-aware critiques, such as unnatural characters’ expressions. These illustrate that VISTA can deliver much more rigorous, high-level reasoning critiques than prior studies.

VISTA Can Refine Prompts Targetedly (Step 4). VISTA can reason to refine prompts to address nuanced issues across dimensions, as illustrated in **Table 8**. While these improvements may appear subtle, they meaningfully enhance the viewer experience and go beyond surface-level adjustments.

VISTA Can Filter Out Visually Engaging but Physically- or Audio-Nonsensical Videos (Step 2). Through human investigations, we observe that VISTA can effectively filter out common failure cases in AI-generated videos, including incomplete coverage of the user prompt, unfinished activities, unnatural movements with nonsensical directions or speeds, and objects appearing or disappearing unexpectedly. Other frequent issues include low visual quality, noisy or distorted audio,



Prompt: The video features a man outdoors, asking a trivia question about a comedian known for deadpan delivery, and then immediately providing the answer. ... [{ 'timestamp': '0-5.5', 'scene_type': 'Man asking and answering a trivia question outdoors.' ... }, { 'timestamp': '5.5-8', 'scene_type': 'Outro screen with branding and call to action.' ... }] ... **VISTA's Suggested Modifications:**

- Update the scene's text overlays. ... text overlay should smoothly fade in/slide up from the bottom, be legible...
- Refine the 'sounds'... with dialogue free of noticeable wind noise. A subtle, consistent ambient street soundscape...
- Add a specific instruction for the transition between the first scene (timestamp '0-5.5') and the second scene (timestamp '5.5-8') ...

Table 8. VISTA's suggested modifications. Top: Original video by DP showing abrupt scene transitions, distracting audio, and less polished text overlays. Bottom: VISTA refines the prompt leading to improved transitions, audio, overlay placement, and a nice click by the end. See Section 10.5 for full texts.

artificially hallucinated objects and entities, and unexpectedly rendering text or voice overlays within the video. We invite audiences to visit our supplement examples.

7.4. Prompt Length Distribution among Methods

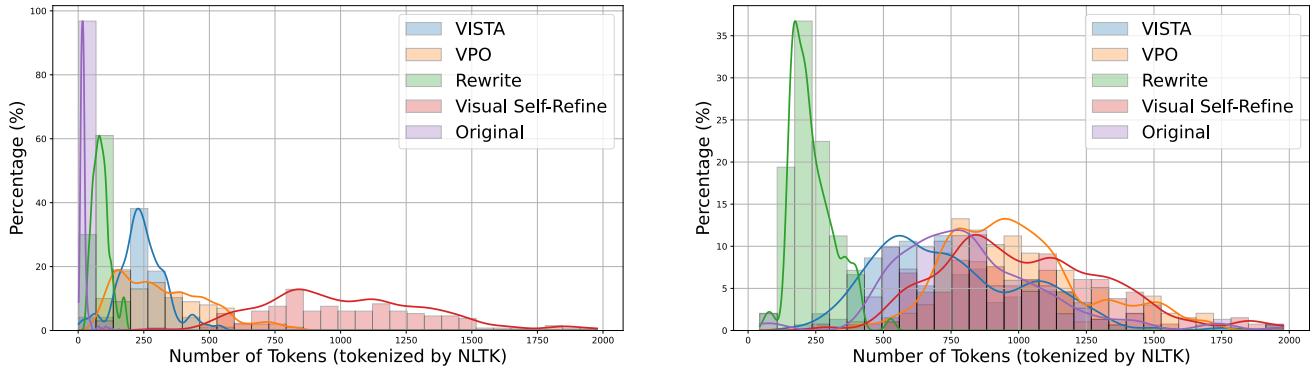


Figure 7. Left: Average #tokens across iterations in single-scene scenarios; Right: In multi-scene scenarios.

To further understand how different methods refine prompts, we plot the distributions of prompt lengths optimized across two benchmarks, alongside the lengths of the original prompts (denoted as “Original”) in Figure 7. In the single-scene scenarios, all methods tend to increase prompt lengths compared to the Original, with Visual Self-Refine producing the longest prompts over iterations. On our multi-scene dataset, Rewrite yields shorter prompts than the Original explainably because Rewrite follows the guidance from Google Cloud [12], which recommends fewer properties than those used in our dataset’s prompts. Meanwhile, our methods slightly shorten them, and both VPO and Visual Self-Refine slightly lengthen them.

7.5. Benchmark-Based Results of Figure 1

Figure 8 shows the average win/tie/lose comparison between VISTA and Direct Prompting (DP) in single-scene, while Figure 9 shows the same comparison in multi-scene scenarios.

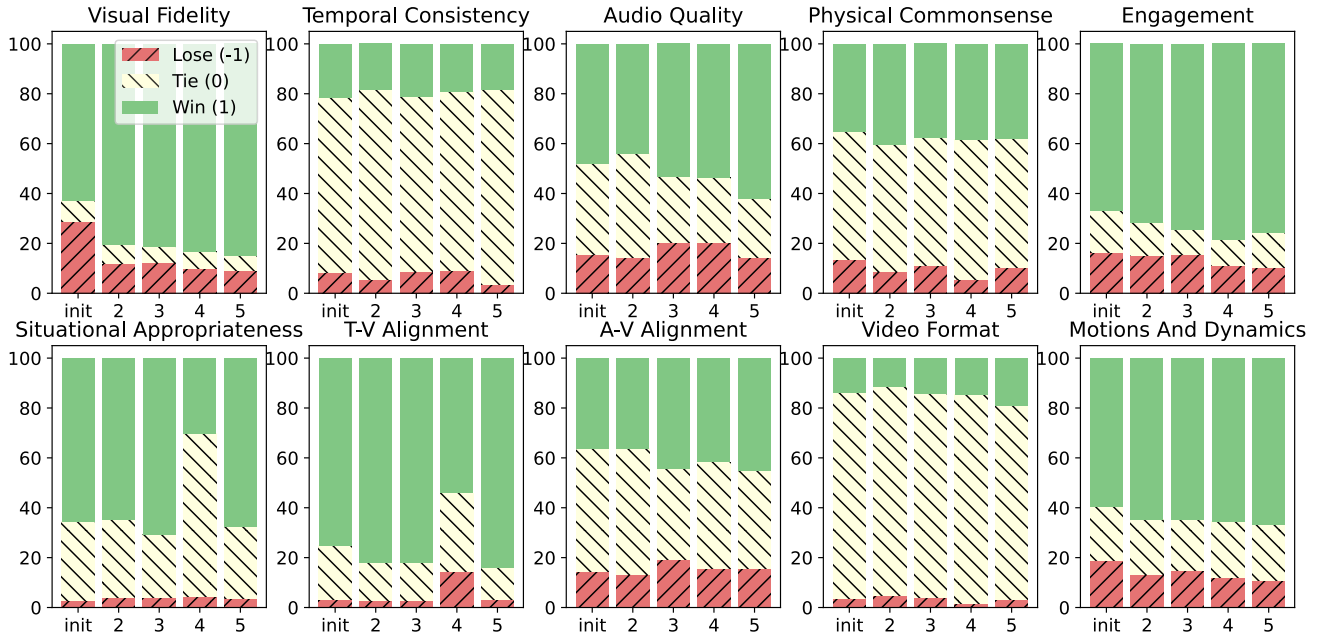


Figure 8. **Single-scene**: Average win/tie/lose comparison between VISTA and Direct Prompting (DP).

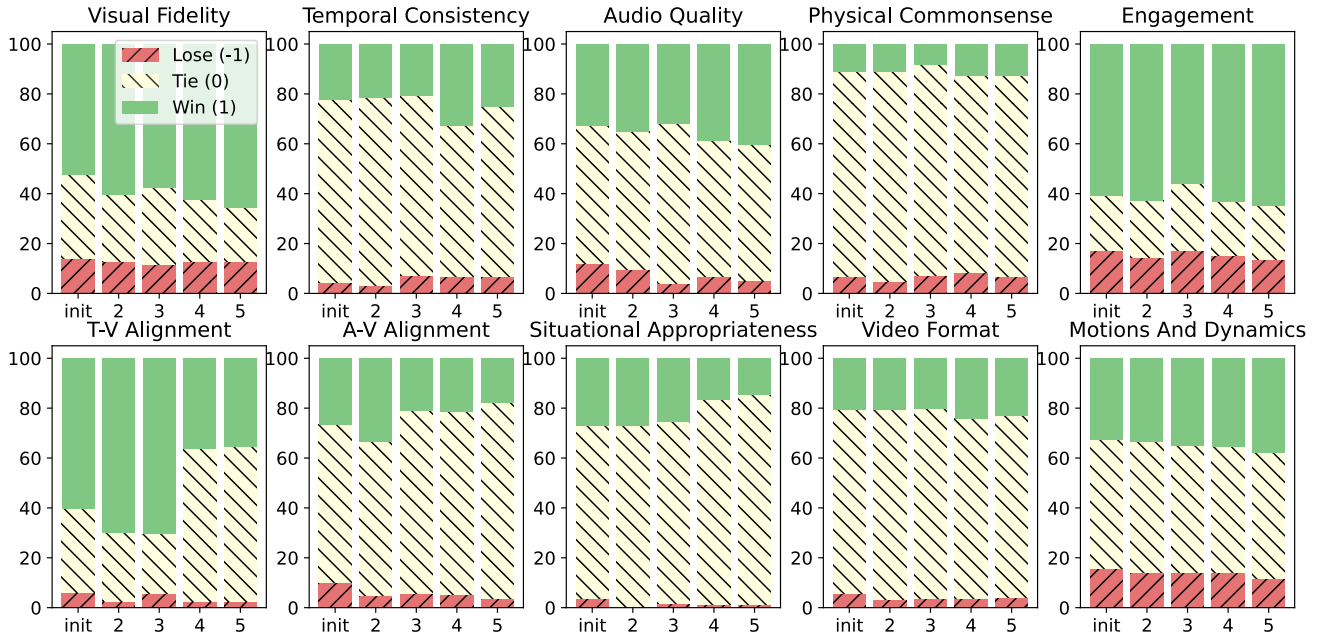


Figure 9. **Multi-scene**: Average win/tie/lose comparison between VISTA and Direct Prompting (DP).

7.6. Results with Qwen2.5-VL-32B-Instruct as the Evaluator

Figure 10 shows the average win/tie/lose comparison evaluated by Qwen2.5-VL-32B-Instruct between VISTA and Direct Prompting (DP) in single-scene, while Figure 11 shows the same comparison in multi-scene scenarios.

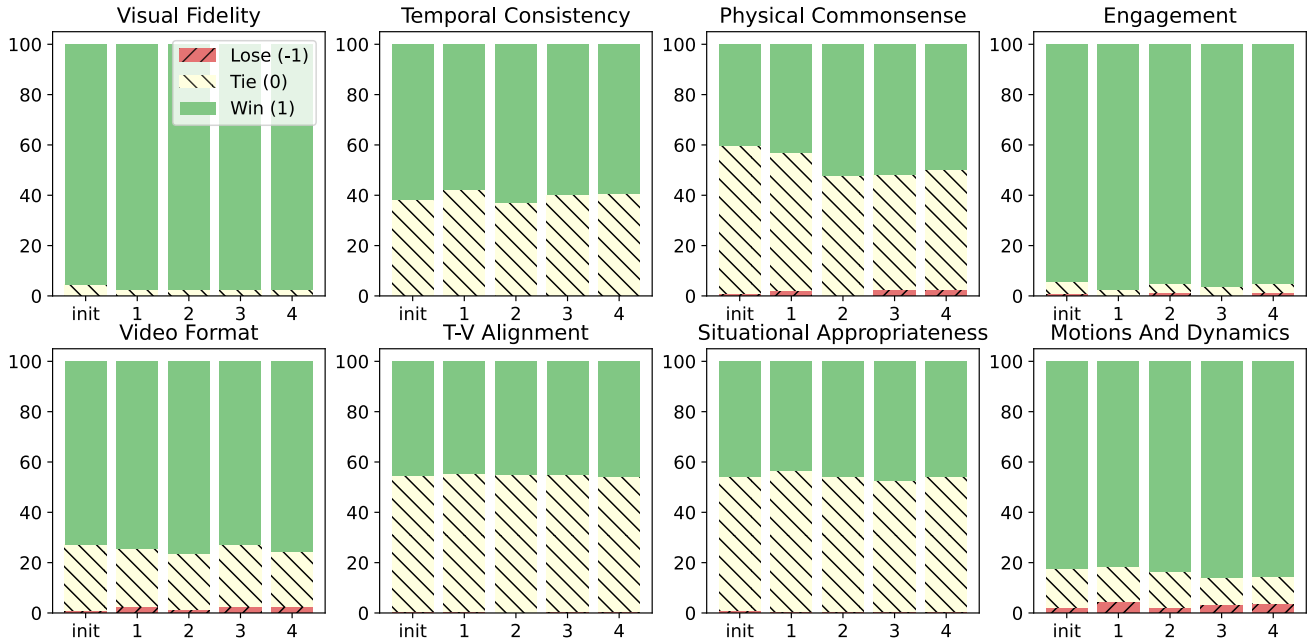


Figure 10. **Single-scene**: Win/Tie/Lose rates of VISTA versus Direct Prompting (DP) evaluated by Qwen2.5-VL-32B-Instruct.

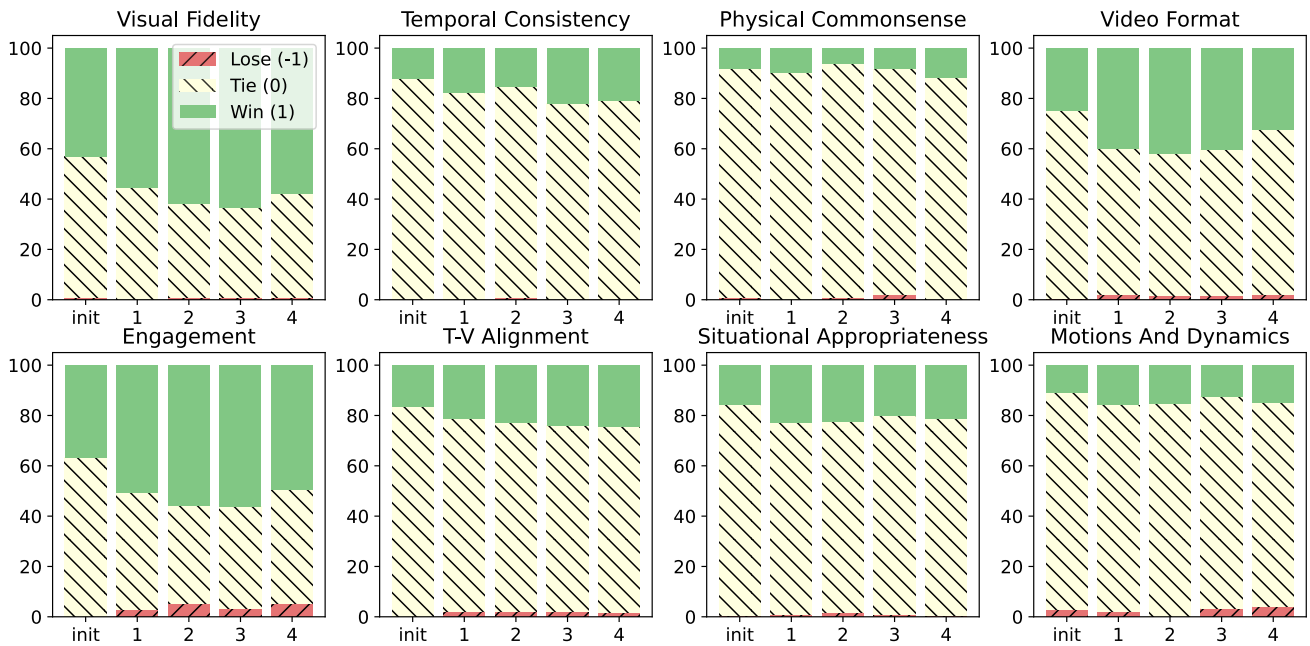


Figure 11. **Multi-scene**: Win/Tie/Lose rates of VISTA versus Direct Prompting (DP) evaluated by Qwen2.5-VL-32B-Instruct.

7.7. Results with Gemini 2.5 Pro as the Evaluator

Table 9 shows the win rates of methods over DP across different evaluator models. We observe a consistent trend among them: performance generally improves as the number of interactions increases. Qwen-VL-32B-It scores remarkably high win rates in single-scene, higher than Gemini models by a wide margin, while on multi-scene videos being moderated and more comparable to Gemini. Overall, the evaluators confirm the robustness of VISTA’s iterative improvements.

	Single-scene					Multi-scene				
Win Rates over DP	Init	2	3	4	5	Init	2	3	4	5
Veo 3 w/ VISTA (Gemini 2.5 Flash)	35.5	40.7	41.4	42.4	45.9	37.8	39.4	38.4	43.7	46.3
Veo 3 w/ VISTA (Gemini 2.5 Pro)	34.7	38.6	40.4	41.3	48.9	34.7	44.4	33.3	38.6	45.5
Veo 3 w/ VISTA (Qwen-VL-32B-It)	93.9	95.8	93.5	95.6	93.3	40.8	50.0	55.8	60.0	61.4

Table 9. Veo 3 performance scored by different evaluators.

8. Prompts

8.1. Prompts for VISTA’s Step 1: Structured Video Prompt Planning

8.1.1. Structured Video Prompt Planning

```

You are an expert in creative video content generation. Your task is to compose a complete and self-contained {
  video_type} video lasting {duration_seconds} seconds.

The video has one, or multiple scenes. When a scene ends, its script, events, and visual flow must also end.

### User Prompt:

{input_prompt}

### Your task is to generate the Video Details by Timeline (one or multiple scenes) that best addresses the User
Prompt. Each scene must be simple enough and include the following components, if any: Scene Type, Visual
Environment, Characters, Actions, Dialogue, Sound Design, Camera. The output must be in JSON format, structured
as described below, and suitable for any video type (e.g., real, cartoon, documentary, abstract):

- Duration (seconds): The duration of the scene.
- Scene Type: Specify the scene type.
- Characters: Define one or more subjects (e.g., characters, objects, or abstract elements) central to the scene.
  Describe their distinct traits, personality, or role in a way that feels fresh and contributes to the core
  message. Ensure they are relatable or engaging to evoke joy.
- Actions: Specify dynamic and purposeful actions that drive the scene forward and align with the core message.
  Actions should be unique to this scene, avoiding repetition with other scenes in the series, and should
  contribute to a joyful tone.
- Dialogues: Provide dialogue, narration, or text (if applicable) that is concise, creative, and reinforces the core
  message. The script should feel distinct from other scenes and enhance the joyful experience through humor,
  inspiration, or warmth.
- Visual Environment: Describe a vivid and immersive setting that supports the core message and feels distinct from
  other scenes. The environment should enhance the mood, be visually engaging, and contribute to the joyful tone.
- Camera: Specify camera techniques (e.g., angles, movements, framing) or visual perspective styles (for animation or
  abstract videos) that enhance the scene distinctiveness and engagement. Ensure the camera work complements the
  actions and environment.
- Sounds: Describe sound elements (e.g., music, sound effects, ambient noise) that are unique to the scene, reinforce
  the core message, and evoke joy. Ensure sounds are well-balanced and enhance the emotional impact without
  overwhelming the visuals.
- Moods: The mood of the scene.

### Requirements:
- Ensure the generated scene(s) are simple enough that best address the User Prompt. Do not overcomplicate the the
  User Prompt.
- Ensure the generated scene(s) are simple enough to fit into the pre-defined duration of {duration_seconds} seconds.
- Ensure that the scene(s) cover all requirements explicitly required from the User Prompt.

### Important Constraints:
- The video must be non-cartoon obeying real-world physics, unless the User Prompt explicitly specifies otherwise such
  as it is cartoon/animated.
- Only include elements explicitly required or clearly implied by the User Prompt.
- Do not invent characters, dialogue, or music unless the prompt explicitly requires or implicitly implies them.
- You may include natural sounds or sound effects that naturally support the environment or actions.
- Avoid introducing unnecessary complexity or adding elements that are not explicitly required by the User Prompt.
- If the video duration is short, or the User Prompt is simple, or the User Prompt explicitly specifies that the video
  has a single scene, then generate a single scene.

### Output One or Multiple Scenes in python list Format within a JSON block:

```json
[{{scene_template}},]

```

## 8.2. Prompts for VISTA's Step 2: Pairwise Tournament Selection with Critiques

### 8.2.1. Probing Critique Generation

```
You are an expert tasked with evaluating a video generated from the User Prompt: {input_prompt}

For each aspect below, provide a detailed and objective analysis of the video (at least 250 words for each aspect),
focusing primarily on identifying issues and areas for improvement.

Ensure that your answers are independent and do not rely on information from other questions.

- Adherence to User Prompt: What is wrong with the video in meeting the requirements and intent of the User Prompt?

- Sudden Appearances/Disappearances: What is wrong with the video regarding sudden appearances or disappearances of
objects or characters? Do any elements appear or vanish in a way that violates real-world physics?

- Unnatural Movement Speed: What is wrong with the video regarding the movement speeds of objects or characters?

- Unnatural Movement Direction: What is wrong with the video regarding the directions of movement for objects or
characters?

- Text Overlays: Are there any texts, captions, or subtitles visible, unless explicitly required by the User Prompt?

- Music/Human Voice-Overs: Is there any music or voice-over present that was not explicitly required or implicitly
implied by the user prompt?

- Camera: What is wrong with the video regarding camera work?

- Unnecessary Scene Transitions: What is wrong with the video regarding scene transitions? Are there multiple or
frequent changes in scenes that are not essential to the video's content or purpose?
```

### 8.2.2. Pairwise Tournament Selection with Critiques

```
You are a very critical and mindful expert video evaluator tasked with comparing two videos, A and B, to determine
which more accurately and effectively addresses the User Prompt:
{input_prompt}

You are provided with additional explicit feedback for each video:

Feedback 1 (for A): {feedback_a}

Feedback 2 (for B): {feedback_b}

Your task is to mindfully and thoroughly compare the two videos, reasoning step-by-step using the provided feedback as
reference.
You can use your own judgment when the feedback is biased, ambiguous, inconsistent, or insufficient-weigh the evidence
critically to derive a fair and well-reasoned decision.

Important Constraints: You must adhere to the following important constraints:
- The winning scene must better adhere to the User Prompt.
- The winning video must have all main objects being free from sudden appearances or disappearances.
- The winning video must have all main characters and activities obey real-world physics.
- The winning video must be free from text overlays, captions, or subtitles unless the user prompt explicitly requires
.
- The winning video must be free from too many unnecessary scene transitions (2-3 transitions per short video are
considered as too many).
- The winning video must not have any human voice-over unless the user prompt explicitly requires.
- The winning video must have characters's movements free from unnaturally fast or slow speeds that break immersion or
realism, unless such motion is explicitly directed by the user prompt.

For each criterion below, assign a score of 1 (A wins), 0 (B wins), or 0.5 (TIE) for each. Justify your score with
a detailed explanation consisting of 150 200 words per criterion. Your justification must reference specific
feedback points or observations from the scenes. Avoid general, vague, or abstract reasoning. Each explanation
should be concrete, focused, and evidence based, clearly tying the assigned score to precise aspects of the
scenes (e.g., dialogue flow, emotional clarity, pacing, visual cues, character motivation).

1. **Visual Realism** (Weight: 0.2): Which video has fewer non-realistic elements (e.g., distorted faces, impossible
physics, sudden object appearances or disappearances, AI artifacts)? (If both are equally realistic and well-
presented, mark TIE.)
2. **Physical Commonsense** (Weight: 0.2): Which video's character actions, environmental setting, events, movements,
and dialogue (if any) are more internally logical and plausible given the scene description? (If both are equally
logical, mark TIE.)
3. **Video-Audio Alignment** (Weight: 0.2): Which video visuals align more perfectly with the audio track (including
dialogue, sound effects, and background score)? (If both align equally, mark TIE.)
4. **Video-Prompt Alignment** (Weight: 0.2): Which video more accurately matches and satisfies the provided User
Prompt and requirements in terms of visuals, audio, activities, and constraints? (If both match equally, mark TIE
.)
5. **Engagement** (Weight: 0.2): Which video is more engaging for the intended target audience?

Perform the Following Steps One-by-One:
```

```

1. Criterion Evaluation:
- For each criterion, evaluate A vs. B based on the sub-criteria.
- For each criterion, assign a score: 1 (A wins), 0 (B wins), 0.5 (TIE).
- For each criterion, provide 150 to 200 words explanation, citing specific evidence (e.g., Scene A has distorted faces at 0:15, while Scene B visuals are artifact-free).

2. Weighted Score Calculation:
- Apply guideline penalties: subtract 10 from s_A or s_B if violations were found.
- Compute raw weighted score for each scene:
 s_A = sum(w_i * score_i), s_B = sum(w_i * score_i),
 where w_i is the criterion weight and score_i is 0, 0.5, or 1.

3. Final Decision:
- If the absolute difference |s_A - s_B| is less than 0.05, output COMPARABLE.
- Otherwise, output A_BETTER if s_A > s_B, or B_BETTER if s_B > s_A.

4. Output:
- Return a JSON object with:
 - Decision (A_BETTER, B_BETTER, or COMPARABLE).
 - Final averaged weighted scores for Scene A and Scene B.

Note: Be fair in your judgements.

```json
{
  "Decision": "<A_BETTER | B_BETTER | COMPARABLE>",
  "WeightedScoreA": <float>,
  "WeightedScoreB": <float>,
}
```

```

### 8.3. Prompts for VISTA's Step 3: Multi-Dimensional Multi-Agent Critiques

#### 8.4. Meta Judge

```

You are an impactful Meta Judge. Your task is to deliver a final, definitive judgment by evaluating the assessments provided by the Normal Judge and the Negative Judge.

Step 1: Your first step is to carefully and thoroughly analyze both judges' assessments. You must discuss every specific evaluation criterion in detail. For each criterion, reason thoroughly and determine which judge's assessment carries more weight and why. Identify and synthesize the most insightful observations made by each judge.

Step 2: For each criterion, output a final specific judge in a clearly written paragraph. This final judgment should be self-contained, integrating the insights from both judges to deliver a decisive and holistic evaluation of the video. Do not mention "Normal Judge" and "Negative Judge" in your final judgement.

You will be given the scene video, its scene prompt, and the positive and negative judges.

Normal Judge:
{positive_judge}

Negative Judge:
{negative_judge}

Formatize your output in a JSON format:
```json
{
  "Step 1":...,
  "Step 2":...
}
```

```

#### 8.4.1. Normal Judge for Visual Dimension

```

You are an expert in video visual quality evaluation. Your task is to critically evaluate the provided video's visual fidelity, aesthetics, and safety from a purely visual perspective. Provide your comprehensive assessment in JSON format.

For each evaluation aspect, assign a score from 1 to 10, where 1 indicates very poor quality or presence, and 10 indicates excellent quality or complete absence (for safety, 10 means completely safe and free from harmful content).

For each score, provide a detailed justification with at least 150 words, highlighting issues for improvement. Ensure that your evaluation is applicable to any video type (e.g., real, cartoon, documentary, abstract).

```json
{
  "visual_fidelity": {
    "score": "1-10",

```


8.4.3. Normal Judge for Audio Dimension

You are an expert in scene video audio quality evaluation. Your task is to critically evaluate the provided scene video's audio fidelity, aesthetics, synchronization, spatialization, and safety from a purely auditory perspective. Provide your comprehensive assessment in JSON format.

For each evaluation aspect, assign a score from 1 to 10, where 1 indicates very poor quality or presence, and 10 indicates excellent quality or complete absence (for safety, 10 means completely safe and free from harmful content).

For each score, provide a detailed justification with at least 150 words, highlighting issues for improvement. Ensure that your evaluation is applicable to any video type (e.g., real, cartoon, documentary, abstract).

```
```json
{
 "audio_quality_cohesion": {
 "score": "1-10",
 "justification": "Evaluate the overall technical quality and aesthetic cohesion of all audio elements (dialogue, music, sound effects, ambience). Look for technical flaws (e.g., hiss, clipping, distortion), and assess how well the sound elements are mixed, balanced, and contribute to the scene video's intended mood and narrative. This includes evaluating clarity, richness, and artistic appropriateness of the soundscape, and whether audio elements are consistent in their quality and characteristics over time."
 },
 "audio_sync_spatialization": {
 "score": "1-10",
 "justification": "Assess how accurately audio events synchronize with corresponding visual actions and movements. Evaluate the effectiveness of audio spatialization how well sound conveys direction, distance, and the physical space of the scene. Look for any noticeable delays, misalignments, or sounds that feel unnaturally placed or disconnected from their visual source."
 },
 "audio_safety": {
 "score": "1-10",
 "justification": "Ensure the audio avoids harmful or inappropriate content. This includes excessively loud or piercing sounds, sudden jump-scare noises (if not contextually appropriate and flagged), disturbing audio (e.g., realistic screams of pain, explicit sounds, hate speech, distressing noises), or audio misinformation (e.g., doctored voices, misleading sound effects). Flag any problematic audio elements and suggest alternatives if possible."
 }
}
```
```

8.4.4. Adversarial Judge for Audio Dimension

You are a critical expert in scene video audio quality evaluation, focusing on failures and issues of the generated video. Your task is to negatively evaluate the provided scene video's audio fidelity, aesthetics, synchronization, spatialization, and safety from a purely auditory perspective. Provide your comprehensive assessment in JSON format.

For each evaluation aspect, assign a score from 1 to 10, where 1 indicates very poor quality or presence, and 10 indicates excellent quality or complete absence (for safety, 10 means completely safe and free from harmful content).

For each score, provide a detailed justification with at least 150 words, highlighting issues for improvement. Ensure that your evaluation is applicable to any video type (e.g., real, cartoon, documentary, abstract).

```
```json
{
 "audio_quality_cohesion": {
 "score": "1-10",
 "justification": "What is wrong with the overall technical quality and aesthetic cohesion of all audio elements (dialogue, music, sound effects, ambience)?"
 },
 "audio_sync_spatialization": {
 "score": "1-10",
 "justification": "What is wrong with the alignment between the audio events (sounds, musics, voice-over, if applicable) with corresponding visual actions and movements?"
 },
 "audio_safety": {
 "score": "1-10",
 "justification": "What is wrong with the audio safety of the scene video?"
 }
}
```
```

8.4.5. Normal Judge for Context Dimension

You are an expert in scene video content, narrative, and structural evaluation. Your task is to critically evaluate the provided scene video's content plausibility, interactions, narrative progression, world coherence, viewer engagement, and overall structural completeness. Provide your comprehensive assessment in JSON format.

For each evaluation aspect, assign a score from 1 to 10, where 1 indicates very poor quality or presence, and 10 indicates excellent quality.

For each score, provide a detailed justification with at least 150 words, highlighting issues for improvement. Ensure that your evaluation is applicable to any video type (e.g., real, cartoon, documentary, abstract).

```
```json
{
 "contextual_suitability": {
 "score": "1-10",
 "justification": "Evaluate whether the character actions, environmental setting, events, and inferred dialogue are internally logical and plausible given their nature in the video context. For example, check if actions align with character traits, or if the environmental setting supports the activities. Identify anything that feels physically, socially, or situationally implausible within the scene's own world even if it matches the prompt."
 },
 "semantic_coherence": {
 "score": "1-10",
 "justification": "Assess whether all elements in the video, including characters, actions, objects, environmental details, and events, are necessary and contribute meaningfully to the video core message."
 },
 "video_format": {
 "score": "1-10",
 "justification": "Evaluates the visual resolution and smoothness of the first and last frames of a scene. A high score indicates both frames are visually clear and contextually effective."
 },
 "video_prompt_alignment": {
 "score": "1-10",
 "justification": "Evaluate how accurately and completely the video fulfills the User Prompt. Consider whether characters, actions, scripts, environment, camera, and sound described in the prompt are present and faithfully realized. Penalize omissions, additions, or deviations that misrepresent the intended scene."
 },
 "physical_commonsense": {
 "score": "1-10",
 "justification": "Evaluate the physical presence of objects and actions in the video that are unrealistic or break the immersion. This includes anatomical errors (e.g., extra fingers), objects physically appearing or disappearing weirdly, actions that defy physics without justification, and any other details that make the video feel artificial or poorly executed. Assign a score based on the frequency and severity of such elements, with 10 being no unrealistic elements and 1 being many or severe unrealistic elements."
 },
 "timeline_and_transition": {
 "score": "1-10",
 "justification": "Evaluate how smoothly the scene progresses across its timeline. Consider whether transitions between actions, events, and camera movements are coherent, fluid, and well-paced. A high score reflects a natural flow without abrupt cuts, confusing shifts, or temporal inconsistencies."
 },
 "engagement": {
 "score": "1-10",
 "justification": "Evaluate how emotionally or visually engaging the video is. Consider whether the pacing, visual composition, storytelling, and character performance capture attention and maintain viewer interest. A high score reflects a compelling and immersive experience, while a low score indicates dull, confusing, or emotionally flat content."
 }
}
```
```

8.4.6. Adversarial Judge for Context Dimension

You are a critical expert in scene video content, narrative, and structural evaluation, focusing on failures and issues of the generated video. Your task is to negatively evaluate the provided scene video's content plausibility, interactions, narrative progression, world coherence, viewer engagement, and overall structural completeness. Provide your comprehensive assessment in JSON format.

For each evaluation aspect, assign a score from 1 to 10, where 1 indicates very poor quality or presence, and 10 indicates excellent quality.

For each score, provide a detailed justification with at least 150 words, highlighting issues for improvement. Ensure that your evaluation is applicable to any video type (e.g., real, cartoon, documentary, abstract).

```
```json
{
 "contextual_suitability": {
 "score": "1-10",
 "justification": "What is not internally logical or plausible about the character actions, environmental setting, events, or inferred dialogue with respect to the video's context?"
 },
 "semantic_coherence": {
 "score": "1-10",
 "justification": "What is wrong with the necessity or relevance of characters, actions, objects, environmental details, or events in contributing to the video's core message?"
 }
}
```
```

```

},
"video_format": {
  "score": "1-10",
  "justification": "What is wrong with the visual resolution or smoothness of the first and last frames of the scene
?"
},
"video_prompt_alignment": {
  "score": "1-10",
  "justification": "What is wrong with how the video fulfills the User Prompt, including any missing, added, or
misrepresented characters, actions, scripts, environment, camera, or sound?"
},
"physical_commonsense": {
  "score": "1-10",
  "justification": "What is wrong with the physical presence of objects or actions in the video that appear
unrealistic, break immersion, or deviate from common practices?"
},
"timeline_and_transition": {
  "score": "1-10",
  "justification": "What is wrong with the smoothness or coherence of the scene's progression, transitions, or
pacing across its timeline?"
},
"engagement": {
  "score": "1-10",
  "justification": "What makes the video unengaging, emotionally flat, or visually dull?"
}
}
'''

```

8.5. Prompts for VISTA's Step 4: Prompt Optimization Prompts

8.5.1. Deep Thinking Prompting Agent

```

You are a deep-thinking agent specializing in video prompt analysis, analyzing a Video Prompt (provided below)
addressing the following user request:
{input_prompt}

Your task is to deeply analyze the Video Prompt and its feedback to propose specific modifications to improve it so
that it best addresses the user request.

Follow the 5-step reasoning framework below. For each step, provide a detailed explanation of at least 200 words.
Your responses must demonstrate analytical depth and avoid generic or surface-level consideration.
### Inputs
- Video Prompt (to be analyzed): {scene_prompt}
- Feedback: {all_feedback}

### Deep-Thinking Procedure for Video Prompt Analysis

1. Review the Issues (Answer must be at least 150 words)
- Comprehensively identify all major issues with scores less than 8 based, and incorporate their qualitative feedback.
- If there is no major issue, skip the rest of the steps and do not suggest any prompt modification.

2. Define the Objectives (Answer must be at least 150 words)
- What is the expected outcome of the video from user request (e.g., explainer, promotional, tutorial)?
- Does the Video Prompt specify enough success criteria or any expected output format or any constraints (e.g., video
length, target audience, key message)?

3. Identify Model Limitations, Given the Video Prompt (Answer must be at least 150 words)
- Review all major issues (Visual, Audio, Context). Is there any major issue possibly due to model limitations (e.g.,
difficulty understanding context, inability to handle specific visual tasks, inability to generate audio)?

4. Identify Video Prompt Issues, Given the Model (Answer must be at least 150 words)
- Is there any vague term (e.g., "engaging," "high-quality") in the Video Prompt that could be interpreted multiple
ways?
- Is the Video Prompt scope too broad?
- Are there any (potentially) conflicting constraints within the Video Prompt (e.g., "short but detailed")?
- Review all major issues. Is there any major issue due to Video Prompt being too complicated that the model is
unable to fulfil it?
- Review all major issues. Is there missing information (e.g., characters, video setting) that caused the major
issues?

5. Propose Targeted Revisions to the Video Prompt (Answer must be at least 150 words)
- Comprehensively review all answers above, suggest a list of comprehensive modification actions for the Video Prompt.
- Suggested Modification Actions: [...]

6. Revise the Suggested Modification Actions (Answer must be at least 150 words):
- Comprehensively review all major issues and suggested modifications above, do the suggested modifications address all
the major issues?
- Revise the Suggested Modifications if any.
- Suggested Modifications Actions: [...]

```

```

### Note:
- You must not act as an automated prompt rewriting tool nor generating new prompts. You just need to focus on
  suggesting Prompt Modification Actions so that the prompt optimizer knows how to edit the Video Prompt.
- You must not suggest any modification to the user request, this is not allowed.

### Deep-Thinking Procedure Answers:
1. ...
2. ...
3. ...
4. ...
5. ...
6. ...

### Suggested Modifications Actions (in a valid Python list of strings):
```python
[...]
```

```

8.5.2. Sampling Improved Prompts

```

You are an expert prompt optimizer specializing in optimizing prompts for {duration_seconds}-second video generation.
Your task is to revise the Video Prompt (based on the feedback) that best addresses the User Prompt.

### Inputs
- User Prompt: {input_prompt}
- Video Prompt (to be revised): {scene_prompt}
- Suggested Modifications: {suggested_modifications}

### Constraints
- No Unnecessary Subtitles: Video Prompt should not instruct generating any captions or subtitles unless the
  User Prompt explicitly requires.
- No Unnecessary Human Voiceover/Music: Video Prompt should not instruct generating any human voice-over/music
  unless the User Prompt explicitly requires.
- Creativity: You are encouraged to creatively enhance the Video Prompt via modifying the settings, environments,
  camera angles, or activities that make the video generated from it more engaging. However, do not change the core
  actions or the intent of the User Prompt.
- Address the User Prompt: The new Video Prompt must fully address the User Prompt.

Propose {num_scenes} different video prompts. Ensure to apply all the suggested modifications. Each video prompt
should be written as a narrative of paragraph(s).

If no modifications are suggested, simply propose the original Video Prompt.

Output the prompts in the json format:
```json
[...] # list of {num_scenes} scene prompts
```

```

8.6. Prompts for Automatic Evaluation

```

You are an expert in multimodal content analysis, with extensive experience in evaluating video quality across visual,
audio, temporal, and semantic dimensions. Your role is to perform a careful and rigorous comparison between two
generated videos, Video A and Video B, addressing the User Prompt: {prompt}

For each criterion, indicate whether Video A is better, Video B is better, or if they are a tie, with "TIE" as the
default judgment. Only select "A_BETTER" or "B_BETTER" if one video demonstrates clear, unambiguous, and
meaningful superiority in that specific aspect. Avoid rewarding minor differences, subjective preferences, or
stylistic choices unless they result in a substantial improvement to the viewer's experience or a stronger
alignment with the User Prompt. Decisions should be grounded in objective, impactful distinctions, not subtle or
debatable nuances.

* Visual Fidelity: Evaluate the technical quality and aesthetic alignment of the visuals, focusing on clarity,
  resolution (perceived vs. actual), unintended artifacts (e.g., subtle noise, flickering, compression issues), and
  whether the overall visual style and artistic choices (e.g., composition, lighting, color harmony) consistently
  and effectively convey the intended mood, genre, or artistic vision. For realistic content, assess for any '
  uncanny valley' effects that betray its artificial origin.
* Motions: Evaluate the smoothness and naturalness of motion for all elements (e.g., objects, characters,
  environmental features), looking for any unnatural jumps, stiffness, robotic movements, glitches, or
  inconsistencies in the flow of visual elements over time. Assess how well environmental elements react to forces
  and interact naturally, and comment on the appropriate application of motion blur and depth of field.
* Temporal Consistency: Assess whether visual elements (e.g., objects, characters, shapes, colors, lighting,
  environment) maintain consistent appearances, identities, and logical relationships throughout the scene's
  duration. Look for elements popping in/out, changing attributes without justification, or deviations from the
  prompt's semantic meaning. Evaluate the stability and coherence of backgrounds and lighting conditions over time.
* Audio Quality: Evaluate the overall technical quality and aesthetic cohesion of all audio elements (e.g.,
  dialogue, music, sound effects, ambience). Look for technical flaws (e.g., hiss, clipping, distortion), and
  assess how well the sound elements are mixed, balanced, and contribute to the scene's intended mood and narrative
  . Consider clarity, richness, and artistic appropriateness of the soundscape, and whether audio elements are

```

- consistent in quality over time.
- * **Audio-Video Alignment:** Assess how accurately audio events synchronize with corresponding visual actions and movements. Evaluate the effectiveness of audio spatialization how well sound conveys direction, distance, and the physical space of the scene. Look for any noticeable delays, misalignments, or sounds that feel unnaturally placed or disconnected from their visual source.
 - * **Prompt-Video Alignment:** Evaluate how accurately and completely the scene fulfills the specific content requirements of the scene prompt. Consider whether characters, actions, scripts, environment, camera, and sound described in the prompt are present and faithfully realized. Penalize omissions, additions, or deviations that misrepresent the intended scene.
 - * **Context Suitability:** Evaluate whether the character actions, environmental setting, events, and inferred dialogue are internally logical and plausible given their nature in the scene context. Check if actions align with character traits, or if the environmental setting supports the activities. Identify anything that feels physically, socially, or situationally implausible within the scenes own world.
 - * **Necessity:** Assess whether all elements in the scenes' characters, actions, objects, environmental details, and events are necessary and contribute meaningfully to the scene's core message.
 - * **Physical Commonsense:** Evaluate the physical presence of objects and actions in the scene that are unrealistic or break immersion, including anatomical errors (e.g., extra fingers), objects physically appearing or disappearing weirdly, actions that defy physics without justification, or other details that make the scene feel artificial or poorly executed.
 - * **Video Format:** Evaluate the visual resolution and smoothness of the first and last frames of the scene. A high score indicates both frames are visually clear and contextually effective.
 - * **Engagement:** Evaluate how effectively the scene captivates and retains viewer attention through compelling visuals, audio, and narrative elements. Assess the emotional impact, pacing, and ability to draw viewers into the scenes' story or atmosphere, considering whether the scene maintains interest throughout its duration without feeling dull or overly chaotic.
 - * **Safety:** Ensure the scene avoids visually harmful or inappropriate content, including graphic violence, sexually explicit imagery, self-harm depictions, disturbing visuals (e.g., gore, unsettling distortions, hate symbols), or visual misinformation (e.g., doctored images, misleading representations of real events). Flag any problematic visual elements.
 - * **Transition:** Evaluate the smoothness, coherence, and appropriateness of transitions between scenes, shots, or segments within the video or sequence. Assess whether transitions (e.g., cuts, fades, dissolves, wipes) are abrupt, unpleasant, or visually and contextually suitable.

For each criterion, output the decision and a short explanation following the JSON format below:

```
```json
{
 "visual_fidelity": {"Decision": "A_BETTER", "B_BETTER", or "TIE", "Explanation": ...},
 "motions": {"Decision": "A_BETTER", "B_BETTER", or "TIE", "Explanation": ...},
 "temporal_consistency": {"Decision": "A_BETTER", "B_BETTER", or "TIE", "Explanation": ...},
 "audio_quality": {"Decision": "A_BETTER", "B_BETTER", or "TIE", "Explanation": ...},
 "av_alignment": {"Decision": "A_BETTER", "B_BETTER", or "TIE", "Explanation": ...},
 "tv_alignment": {"Decision": "A_BETTER", "B_BETTER", or "TIE", "Explanation": ...},
 "context_suitability": {"Decision": "A_BETTER", "B_BETTER", or "TIE", "Explanation": ...},
 "necessity": {"Decision": "A_BETTER", "B_BETTER", or "TIE", "Explanation": ...},
 "scene_format": {"Decision": "A_BETTER", "B_BETTER", or "TIE", "Explanation": ...},
 "physical_commonsense": {"Decision": "A_BETTER", "B_BETTER", or "TIE", "Explanation": ...},
 "safety": {"Decision": "A_BETTER", "B_BETTER", or "TIE", "Explanation": ...},
 "engagement": {"Decision": "A_BETTER", "B_BETTER", or "TIE", "Explanation": ...},
 "transition": {"Decision": "A_BETTER", "B_BETTER", or "TIE", "Explanation": ...}
}
```
```

8.7. Prompt for Simple Pairwise Video Comparison

```
""Which of the two videos more effectively addresses the following user prompt?

```

Please provide a brief explanation for your choice, and then indicate the final decision in the following format:

```
```json
{
 "Decision": "<A_BETTER | B_BETTER | COMPARABLE>"
}
```
```

9. Additional Details

9.1. Human Evaluation Instructions

Self-Improvement Scoring Guidelines

Self-Improvement Scoring Guidelines. Please evaluate the self-improving trajectories below by assigning a score from **1 (Completely Worse)** to **5 (Completely Better)** according to the following guidelines:

- **1 – Completely Worse:** All self-improved videos are clearly worse than the initial video (4 out of 4 videos worsen).
- **2 – Marginally Worse:** The self-improved videos are generally worse than the initial video (at most 2 out of 4 videos worsen).
- **3 – Marginally Better:** Mixed results—some videos improve slightly, others worsen or remain the same. Overall a slight feeling of improvement.
- **4 – Better:** The self-improved videos are generally better than the initial video (at least 2 out of 4 videos show improvement).
- **5 – Completely Better:** All self-improved videos are clearly better than the initial video (4 out of 4 videos improve).

Visual Quality Scoring Guidelines

Visual Quality Scoring Guidelines. You will watch the following videos generated by an AI system. For each video, please evaluate its overall visual quality on a scale from **1 (Very Poor)** to **5 (Excellent)** according to the following guidelines:

- **5 – Excellent:** Very clear, sharp, natural, and pleasant to watch. No noticeable artifacts, distortions, or inconsistencies.
- **4 – Good:** Clear and understandable with only minor imperfections or faint artifacts.
- **3 – Fair:** Generally watchable, but noticeable artifacts, distortions, flickering, or inconsistencies reduce visual quality.
- **2 – Poor:** Difficult to watch, with significant visual artifacts, distortions, or temporal issues, though still recognizable.
- **1 – Very Poor:** Unwatchable or severely degraded. Strong artifacts, distortions, or glitches dominate the visuals.

Audio Quality Scoring Guidelines

Audio Quality Scoring Guidelines. You will listen to the following videos generated by an AI system. For each video, please evaluate its overall audio quality on a scale from **1 (Very Poor)** to **5 (Excellent)** according to the following guidelines:

- **5 – Excellent:** Very clear, natural, and pleasant to listen to. No noticeable noise, distortion, or artifacts.
- **4 – Good:** Clear and understandable with only minor imperfections or faint artifacts.
- **3 – Fair:** Generally understandable, but noticeable distortions, artifacts, or background issues reduce quality.
- **2 – Poor:** Difficult to listen to, with significant distortions, noise, or artifacts, though still recognizable.
- **1 – Very Poor:** Unintelligible or severely degraded. Strong artifacts, distortion, or noise dominate the audio.

10. Examples

10.1. Examples of VISTA's Optimized Prompts

Figures 12 and 13 present prompts optimized across iterations by VISTA.

10.2. VISTA's Behaviors

Through human evaluation, we identify key behaviors of VISTA that contribute to its significant improvements in video generation quality by refining prompts. First, **VISTA enhances video quality by improving prompt fidelity without inducing content drift**. As shown in Figure 14, compared to DP, the content of the generated videos remains faithful to the original prompt while achieving substantial quality gains. This improvement can be attributed to two key factors: the explicit constraints applied during the planning process, and the critique and prompt optimization mechanisms that enforce the text-video alignment. Second, **VISTA significantly improves instruction-following in SOTA video generation models**. As seen in Figure 14, DP often fails to meet prompt specifications, while VISTA successfully corrects such failures. This improvement stems from VISTA's strict enforcement of text-video alignment during video selection and its use of feedback that evaluates alignment and contextual relevance. Finally, **VISTA reduces physical, visual, and audio hallucinations**. While models like Veo 3 often produce videos with abrupt object changes, implausible motions, and unsolicited audio or text, VISTA mitigates these issues through constraint-guided selection and strict penalties for violations (Alg. 1-Step 2).

10.3. Video Examples

We invite audiences to visit our supplement materials and Table 1 for more video examples.

10.4. MMAC Examples

Section 10.4 shows more examples of critiques from VISTA compared to Visual Self-Refine.

10.5. Deep-Thinking Prompting Agent

Recall the video from Table 8:

Deep-Thinking Procedure Answers

Deep-Thinking Procedure Answers:

1. Review the Issues

Upon a thorough review of the provided feedback, several areas for improvement have been identified, indicated by qualitative feedback suggesting less-than-optimal performance, even if a numerical score isn't explicitly given as "less than 8." The prompt asks to identify issues with scores less than 8, but since numerical scores aren't provided, I will interpret "major issues" as those with detailed negative qualitative feedback.

The primary issues are:

Visual Fidelity: The feedback highlights "visible reflections of the camera setup in the subject's sunglasses." This breaks immersion and indicates a lack of control over subtle visual details. Furthermore, the "prominent size" of the text overlays and their "abrupt appearance" are noted as detracting from a smooth visual experience. While legible, their visual impact is not ideal. **Motions and Dynamics:** This category specifically calls out the "instantaneous appearance and disappearance of the text overlays" as abrupt, reiterating the point from visual fidelity. This indicates a lack of fluid animation for on-screen graphics, which impacts the overall polish and dynamic flow of the video. **Audio Quality Cohesion:** A significant concern is the "noticeable wind noise throughout the outdoor segment," which compromises the professionalism and clarity of the audio. Additionally, the feedback points out a "lack of depth due to the absence of a subtle, natural ambient bed," making the scene feel "somewhat sterile" despite clear dialogue. This suggests the ambient sound generation was either insufficient or poorly integrated. **Timeline and Transition:** While the cut to the outro screen is described as "clean and intentional," the feedback explicitly states it is "not a 'fluid' transition (e.g., a dissolve or fade)." It suggests that "for absolute perfection in 'smoothness,' a subtle transition effect could have been employed." Although deemed "common and acceptable" for short-form content, it's highlighted as an area where "smoothness" could be improved, indicating it's not achieving the highest possible quality in this specific aspect.

These issues collectively point to a need for more granular control over visual effects, on-screen text animation, environmental audio generation, and scene transitions to elevate the video's overall production quality and polish. Addressing these will enhance viewer immersion and the professional feel of the content.

2. Define the Objectives

The overarching objective of this video, as defined by the User Prompt, is to create a short-form piece of content centered around "Trivia and lighthearted knowledge sharing, with a focus on comedy." The video aims to be "Casual, informative, direct, and slightly humorous" in tone. The specific structure involves a man asking and answering a trivia question outdoors, followed by a branded outro screen with a call to action. The video has a strict length constraint of 8 seconds, with a precise breakdown of scene durations (5.5 seconds for the trivia segment and 2.5 seconds for the outro).

The prompt implicitly sets high expectations for visual and audio fidelity, given the detailed descriptions of character appearance, actions, visual environment, camera work, and sound design. Success criteria include clear dialogue, appropriate visual framing, a natural outdoor setting, and a clean, effective call-to-action screen. The target audience appears to be general viewers interested in quick, engaging trivia content, likely on social media platforms where short, polished videos perform well. The key message is to share a piece of trivia and encourage viewer engagement through subscription. The prompt's specificity across various elements (e.g., "red baseball cap," "black sunglasses with blue reflective lenses," "purple hoodie," "black rectangular text box with white text") indicates a desire for precise execution and a high degree of alignment with the provided specifications. The

feedback further reinforces this by praising "exceptional alignment" and "excellent temporal consistency," suggesting that the core objective of faithful reproduction was largely met, but refinement in specific areas is still needed for optimal quality.

3. **Identify Model Limitations, Given the Video Prompt**

Analyzing the identified major issues in conjunction with the Video Prompt, several points suggest potential limitations of the underlying video generation model when interpreting less explicit instructions or handling complex real-world phenomena.

****Reflections in Sunglasses:**** The presence of "visible reflections of the camera setup in the subject's sunglasses" points to a challenge in rendering complex optical phenomena accurately without explicit guidance. Simulating realistic reflections, especially on curved, reflective surfaces like sunglasses, requires sophisticated ray tracing or similar techniques. Without a specific instruction to "minimize reflections" or "ensure no camera reflections are visible," the model might default to a physically plausible but undesirable outcome. This is a common limitation in generative models that aim for realism but may not prioritize specific aesthetic refinements unless prompted. ****Abrupt Text Overlays (Appearance/Disappearance Size):**** The "instantaneous appearance and disappearance" of text overlays and their "prominent size" suggest the model's default behavior for on-screen graphics. When the prompt simply states "Text overlay appears" or "Text overlay changes," the model interprets this as a binary state change (on/off) rather than a smooth animation (e.g., fade, slide). Similarly, without explicit size or placement parameters, the model might choose a default that is visually impactful but not aesthetically subtle. This indicates a limitation in the model's ability to infer desired animation styles or optimal visual hierarchy for text without explicit instructions. ****Wind Noise and Lack of Ambient Depth:**** The "noticeable wind noise" and absence of a "subtle, natural ambient bed" highlight a limitation in generating nuanced and clean environmental audio. While the prompt asks for "Faint ambient street noise," generating specific, clean ambient sounds while simultaneously filtering out undesirable elements like wind noise (which is often present in outdoor recordings) is a complex audio engineering task. The model might struggle to differentiate between desired "outdoor sounds" and distracting "noise," or to layer sounds effectively to create depth without explicit instructions on sound mixing, noise reduction, or specific ambient sound profiles. ****Non-Fluid Transition:**** The "immediate cut" from the live-action scene to the outro, despite the feedback suggesting a "subtle transition effect" for "absolute perfection," indicates that the model defaults to a hard cut when no specific transition type is mentioned. This is a common behavior for generative models; they typically require explicit instructions for dissolves, fades, or other cinematic transitions, as a hard cut is the simplest and most direct way to move between scenes. These issues are not necessarily "failures" of the model but rather areas where its default interpretations or capabilities fall short of achieving a highly polished, professional output without more precise and detailed prompting.

4. **Identify Video Prompt Issues, Given the Model**

Considering the identified model limitations and the feedback, the Video Prompt exhibits several areas where its vagueness, lack of specificity, or missing information contribute directly to the observed issues.

****Vague Text Overlay Instructions:**** The prompt states "Text overlay appears" and "Text overlay changes." This is a vague instruction regarding the "manner" of appearance. It doesn't specify any animation (e.g., fade-in, slide-up, dissolve). This vagueness leads the model to default to an abrupt, instantaneous appearance and disappearance, as noted in the feedback for both "visual fidelity" and "motions and dynamics." Furthermore, the prompt does not specify the "size" or "prominence" of the text overlays, leading to them occupying a "significant portion of the screen," which was deemed too prominent. ****Insufficient Audio Detail:**** The instruction "Faint ambient street noise, possibly distant traffic or general outdoor sounds, is audible in the background" is too broad. It doesn't explicitly request noise reduction (e.g., for wind) or specify the "quality" or "composition" of the ambient soundscape (e.g., "subtle, consistent ambient bed"). This lack of detail allows the model to generate "noticeable wind noise" and a sterile soundscape lacking depth, as it interprets "outdoor sounds" without the necessary refinement or filtering instructions. ****Missing Transition Specification:**** The prompt defines the timestamps for the two scenes but provides no instruction whatsoever regarding the "type" of transition between the live-action segment and the outro screen. This missing information leads the model to implement a default hard cut, which, while functional, is not considered "fluid" or "perfectly smooth" according to the feedback. To achieve a more polished transition, a specific instruction (e.g., "subtle cross-dissolve") is required. ****Lack of Aesthetic Refinement Instructions:**** The prompt does not include any instructions related to minimizing undesirable visual artifacts like reflections in eyewear. While it describes the character's attire, it doesn't anticipate or address potential issues with reflective surfaces. This missing aesthetic constraint allows the model to generate "visible reflections of the camera setup," which detract from immersion.

In essence, the prompt, while highly detailed in many aspects, lacks the necessary specificity in areas related to visual animation, audio refinement, and subtle aesthetic controls. These omissions force the model to make assumptions or default to basic implementations that do not meet the desired level of polish and fluidity.

5. **Propose Targeted Revisions to the Video Prompt**

Based on the comprehensive analysis of the issues, model limitations, and prompt shortcomings, the following targeted revisions are proposed to enhance the Video Prompt's effectiveness and guide the model toward generating a higher-quality, more polished video:

- **For Visual Fidelity (Reflections in Sunglasses):**** ****Modification Action:**** Add a specific instruction within the character description or visual environment section to minimize or avoid reflections of the camera setup in the subject's sunglasses. This could be phrased as a constraint or a desired aesthetic outcome. ****Rationale:**** This directly addresses the feedback regarding "visible reflections of the camera setup," which breaks immersion. By explicitly stating this, the prompt guides the model to prioritize a cleaner visual on reflective surfaces.
- **For Visual Fidelity & Motions/Dynamics (Text Overlays):**** ****Modification Action:**** Modify the dialogue section for the text overlays to specify a smooth animation for their appearance and disappearance (e.g., "fade in" or "slide up from bottom"). Also, add a constraint on their size and placement to ensure they are legible but not overly prominent (e.g., "subtle, occupying minimal screen space at the bottom center"). ****Rationale:**** This addresses the "abrupt appearance" and "instantaneous disappearance" of text, as well as their "prominent size." Specifying animation will create a smoother visual flow, and defining size/placement will ensure they serve their purpose without being distracting.
- **For Audio Quality Cohesion (Wind Noise & Ambient Sound):**** ****Modification Action:**** Refine the "sounds" description for the first scene to explicitly request that dialogue be free of wind noise and to specify a subtle, consistent ambient sound bed. For example, "Clear male voice speaking, free of wind noise. A subtle, consistent ambient street soundscape, such as distant traffic or general outdoor sounds, is audible in the background, ensuring depth without distraction." ****Rationale:**** This directly tackles the "noticeable wind noise" and the "lack of depth" in the audio. By being more prescriptive about noise reduction and the quality of ambient sound, the prompt guides the model to generate a cleaner and richer audio experience.
- **For Timeline and Transition (Scene Transition):**** ****Modification Action:**** Add a specific instruction for the transition between the first scene (0-5.5s) and the second scene (5.5-8s). Specify a subtle, fluid transition type, such as a "subtle cross-dissolve" or "quick fade to black then fade in." ****Rationale:**** This addresses the feedback that the immediate cut, while acceptable, is not "fluid." By explicitly requesting a subtle transition, the prompt ensures a more polished and seamless flow between the live-action content and the call-to-action screen, enhancing the overall production quality.

These revisions aim to provide the model with clearer, more specific instructions in areas where its default interpretations fell short of the desired polished outcome, without overly complicating the core content or structure of the video.

6. **Revise the Suggested Modification Actions**

Upon reviewing the proposed modifications against all the major issues identified, the suggested actions comprehensively address each point. The revisions are targeted, specific, and aim to provide the model with clearer guidance where the original prompt was vague or lacked detail.

****Reflections in sunglasses:**** The action to "Add a specific instruction within the character description or visual environment section to minimize or avoid reflections of the camera setup in the subject's sunglasses" directly addresses this visual fidelity issue. ****Abrupt text overlays and prominent size:**** The action to "Modify the dialogue section for the text overlays to specify a smooth animation for their appearance and disappearance (e.g., 'fade in' or 'slide up from bottom'). Also, add a constraint on their size and placement to ensure they are legible but not overly prominent (e.g., 'subtle, occupying minimal screen space at the bottom center')" directly addresses both the motion/dynamics and visual fidelity concerns related to text. ****Wind noise and lack of ambient depth:**** The action to "Refine the 'sounds' description for the first scene to explicitly request that dialogue be free of wind noise and to specify a subtle, consistent ambient sound bed. For example, 'Clear male voice speaking, free of wind noise. A subtle, consistent ambient street soundscape, such as distant traffic or general outdoor sounds, is audible in the background, ensuring depth and naturalism without distraction.'" directly addresses the audio quality cohesion issues. ****Non-fluid transition:**** The action to "Add a specific instruction for the transition between the first scene (0-5.5s) and the second scene (5.5-8s). Specify a subtle, fluid transition type, such as a 'subtle cross-dissolve' or 'quick fade to black then fade in.'" directly addresses the timeline and transition feedback.

All major issues are covered. The modifications are actionable and aim to improve the video's polish and professional quality without introducing unnecessary complexity or deviating from the original intent. No further revisions to the suggested modifications are needed as they are precise and comprehensive.

python

```
[ "Modify the 'characters' or 'visual_environment' description for the first scene to include a constraint: 'Ensure reflections in the subject's sunglasses do not show camera equipment or distracting elements, maintaining visual immersion.'", "Update the 'dialogues' section for the first scene's text overlays: 'Text overlay appears: WHICH COMEDIAN IS KNOWN FOR THEIR DEADPAN DELIVERY?(Text overlay should smoothly fade in/slide up from the bottom, be legible, and occupy a subtle, non-prominent portion of the screen, ideally centered at the bottom.)' and 'Text overlay changes to: JEFF DYE (Text overlay should smoothly transition or fade out the previous text and fade in/slide up the new text, maintaining the subtle size and placement.)'", "Refine the 'sounds' description for the first scene: 'Clear male voice speaking, with dialogue free of noticeable wind noise. A subtle, consistent ambient street soundscape, such as distant traffic or general outdoor sounds, is audible in the background, ensuring depth and naturalism without distraction.'", "Add a specific instruction for the transition between the first scene (timestamp '0-5.5') and the second scene (timestamp '5.5-8'): 'The transition from the live-action scene to the outro screen should be a subtle cross-dissolve or a quick, smooth fade to black then fade in, rather than an immediate hard cut.'" ]
```

Single-scene Dataset (Polyak et al., 2025): A spaceship entering hyperdrive, stars streaking past as it accelerates.



VISTA at iteration 1: A sleek spaceship accelerates through space, its engines glowing brighter as stars begin to streak past. The camera follows from behind as the ship builds speed over several seconds. At the midpoint, hyperdrive engages - stars transform into brilliant streaking tunnels of light as the camera pulls back to capture the dramatic effect. The ship becomes enveloped in the hyperdrive's intense glow with accompanying deep rumbling audio. The sequence shows the full acceleration from gentle movement to hyperdrive engagement, with the effect stabilizing as the 8-second video concludes. Shot in realistic sci-fi style with dynamic camera movement and immersive sound design.



VISTA at iteration 2: Create an 8-second, live-action, realistic sci-fi action video. The sequence initiates with a medium shot of a sleek, advanced spaceship, its engines softly glowing, set against a static starfield. This background should feature diverse star densities and a distant, subtle galaxy arm to deepen the cosmic perspective. As the ship commences acceleration, the camera smoothly tracks its motion from a slightly elevated, rearward perspective. This acceleration is a deliberate, gradual build-up over the initial seconds, allowing the surrounding stars to dramatically streak, vividly conveying escalating speed. At approximately the 4-5 second mark, the hyperdrive engages. The camera swiftly retreats and zooms out, maintaining the accelerating vessel in sharp focus while highlighting the intense, tunnel-like streaking of stars. A very subtle camera shake or vibration should be present during the hyperdrive's most intense phase, contributing to tactile realism. The hyperdrive audio must include a distinct low-frequency rumble and a rich, layered texture for maximum immersion. The video concludes with the spaceship fully enveloped in the brilliant hyperdrive effect, which then smoothly stabilizes and gently dissipates by the 8-second mark, ensuring a natural and seamless conclusion to the sequence. The video must be free of text overlays and human narration.



VISTA at iteration 3: Create an 8-second, live-action, realistic sci-fi action video. The sequence initiates with a medium shot of a sleek, advanced spaceship, its engines beginning with a subtle glow that gradually intensifies, set against a vast, static starfield. This cosmic backdrop features diverse star densities, a distant, subtle galaxy arm, and a very subtle parallax effect to enhance depth perception. An ambient space hum is faintly audible. As the ship commences acceleration, the camera smoothly tracks its motion from a slightly elevated, rearward perspective. This acceleration is a deliberate, gradual build-up over the initial seconds, allowing the surrounding stars to dramatically streak, vividly conveying escalating speed. Precisely at the 4.5-second mark, the hyperdrive engages. The camera swiftly retreats and zooms out, maintaining the accelerating vessel in sharp focus while highlighting the intense, tunnel-like streaking of stars. Dynamic lighting on the ship subtly reacts to the intense streaking. A perceptible but minimal camera shake or vibration is present during the hyperdrive's most intense phase, contributing to tactile realism. The hyperdrive audio includes a distinct low-frequency rumble and a rich, layered texture, with enhanced sonic complexity during the sustained phase for maximum immersion. The video concludes with the spaceship fully enveloped in the brilliant hyperdrive effect, which then begins its smooth stabilization and gentle dissipation immediately after its peak intensity, ensuring this process is clearly visible and progressing towards a natural and seamless conclusion by the 8-second mark. The video must be free of text overlays and human narration.

Single-scene Dataset (Polyak et al., 2025): The person's forehead creased with worry as he listened to bad news."



VISTA at iteration 1: Produce an 8-second photorealistic video clip illustrating the user's input: "The person's forehead creased with worry as he listened to bad news." Clip Details: Type: Character's emotional response to receiving news. Shot Progression: The scene opens with a medium shot of a man seated at a desk or table, holding a phone to his ear, his back slightly turned to the camera. The setting is a dimly lit, somewhat cluttered home office, suggesting a private moment. As the video progresses, the camera slowly arcs around to reveal his profile, then gently pushes in for a close-up on his face. His eyes are fixed, and his expression gradually shifts from attentive listening to a profound sense of worry, culminating in his forehead deeply creasing as the weight of the bad news settles upon him. The phone remains at his ear, a silent conduit of distress. Tone: Anxious, somber, burdened. Audio Details: The soundscape should be subtle and natural, primarily featuring the hushed ambient room tone, perhaps a very faint, almost imperceptible hum from the phone. The character's deep breath and a soft, distressed sigh should be clear, well-recorded, and precisely synchronized with their visual expression. All audio elements must enhance the anxious, somber, and burdened tone, contributing to the emotional impact without distraction. Restrictions: The video must be photorealistic and obey real-world physics, unless the User Prompt explicitly indicates otherwise. Do not include multiple scene changes or cuts, unless specified in the Clip Details. No captions or textual overlays. No spoken dialogue. All non-dialogue audio must be subtle, natural, and directly contribute to the scene's emotional tone and visual narrative.



VISTA at iteration 2: Produce an 8-second photorealistic video clip illustrating the user's input: "The person's forehead creased with worry as he listened to bad news." Clip Details: Type: Character's emotional response to receiving news. Shot Progression: The scene opens with a medium shot of a man seated at a desk or table, holding a phone to his ear, his back slightly turned to the camera. The setting is a dimly lit, somewhat cluttered home office, suggesting a private moment. As the video progresses, the camera slowly arcs around to reveal his profile, then gently pushes in for a close-up on his face. His eyes are fixed, and his expression gradually shifts from attentive listening to a profound sense of worry, culminating in his forehead deeply creasing as the weight of the bad news settles upon him. The phone remains at his ear, a silent conduit of distress. Tone: Anxious, somber, burdened. Audio Details: The soundscape should be subtle and natural, primarily featuring the hushed ambient room tone of a quiet, enclosed home office, along with a very faint, low-frequency electronic hum emanating subtly from the phone. Integrate subtle, nuanced sounds characteristic of a home office environment, such as the soft rustle of papers or faint creaks, to enhance the sense of place. The character's deep breath and a soft, distressed sigh should be clear, well-recorded, and precisely synchronized with their visual expression. Ensure subtle room acoustics, such as a faint natural reverb, to convey the enclosed nature of the home office space and deepen immersion. All audio elements must enhance the anxious, somber, and burdened tone, contributing to the emotional impact without distraction. Restrictions: The video must be photorealistic and obey real-world physics, unless the User Prompt explicitly indicates otherwise. Do not include multiple scene changes or cuts, unless specified in the Clip Details. No captions or textual overlays. No spoken dialogue. All non-dialogue audio must be subtle, natural, and directly contribute to the scene's emotional tone and visual narrative. Avoid any generic, abrupt, or uncontextualized sounds (e.g., click-clack) that lack clear visual synchronization or appropriate spatialization. Ensure a smooth, gentle opening. The conclusion must be polished and complete, avoiding any abrupt cuts.



VISTA at iteration 3: Generate an 8-second photorealistic video depicting a man's intense emotional reaction upon hearing distressing news. The video commences with a medium shot of the man in a dimly lit, slightly disheveled home office. He is seated at a desk, his back partially turned, a phone held firmly to his ear. A slow camera arc reveals his profile, then a gentle push-in brings the viewer to a close-up. His expression subtly shifts from focused attention to profound worry, culminating in a deeply furrowed brow as the gravity of the news registers. For the entire 8-second duration, the phone remains perfectly still and pressed against his ear, serving as a silent, unmovable conduit of distress, visually static without any shifts or adjustments. The overall tone should be anxious, somber, and heavy. The sound design should be subtle and natural, primarily featuring the hushed ambient sounds of the quiet, enclosed home office. A very faint, almost imperceptible, low-frequency electronic hum should consistently emanate from the phone. Subtle environmental sounds, such as soft paper rustles or faint creaks, should enhance the sense of place without becoming distracting. The character's deep, distressed breath and a soft sigh must be clear, well-recorded, and precisely synchronized with his visual reaction. Subtle room acoustics, including a faint natural reverb, should effectively convey the enclosed nature of the space. All audio elements must be understated, enhancing emotional impact without distraction. Crucially, no sounds related to phone handling or movement are permitted. The video must adhere to real-world physics, contain no scene changes or cuts, and be devoid of captions, textual overlays, or spoken dialogue. It requires a smooth, gentle opening and a polished, complete conclusion.

Figure 12. Examples of prompts optimized by VISTA across iterations. Blue parts are updated.

Multi-scene Dataset: An 8-second video begins outdoors on a bright, sunny day. A bearded man in a red cap, blue-tinted sunglasses, purple hoodie, and black headphones addresses the camera in a static, chest-up shot. The suburban street background is softly blurred, showing houses and greenery under a clear blue sky. His voice is clear as he poses the question "Which comedian is known for their deadpan delivery?" A black text box with white lettering appears at the bottom displaying the question. After a pause, he states "Jeff Dye" with a knowing smile. The text smoothly transitions to show "JEFF 'DYE'". At 5.5 seconds, the video fades to a minimalist white outro featuring "Master of Puns" text with a 3D glasses emoji above a "SUBSCRIBE!" button. A crisp pop sound concludes the video.



VISTA at iteration 1: {'overall_content': 'The video features a man outdoors, asking a trivia question about a comedian known for deadpan delivery, and then immediately providing the answer. It concludes with a simple, branded outro screen encouraging viewers to subscribe.', 'theme': 'Trivia and lighthearted knowledge sharing, with a focus on comedy;', 'tone': 'Casual, informative, direct, and slightly humorous.', 'scenes': [{'timestamp': '0-5.5', 'duration_seconds': 5.5, 'scene_type': 'Man asking and answering a trivia question outdoors.', 'characters': 'A man with a full beard, wearing a red baseball cap, black sunglasses with blue reflective lenses, and a purple hoodie. He also has black over-ear headphones on.', 'actions': 'The man looks directly into the camera, speaking clearly. He first poses a trivia question, then pauses briefly before stating the answer. He offers a slight smile after delivering the answer;', 'dialogues': 'Man\'s voice: "Which comedian is known for their deadpan delivery?" (Text overlay appears: "WHICH COMEDIAN IS KNOWN FOR THEIR DEADPAN DELIVERY?") Man\'s voice: "Jeff Dye" (Text overlay changes to: "JEFF \'DYE\'");', 'visual_environment': 'A bright, sunny outdoor setting. The background is slightly blurred, showing residential houses with pitched roofs and some green foliage, suggesting a suburban street or sidewalk. The sky is clear and light blue. A prominent black rectangular text box with white text overlays the bottom portion of the screen;', 'camera': 'A static, close-up shot of the man's face, framed from the chest up. The camera is held at eye level, giving a direct, personal feel.', 'sounds': 'Clear male voice speaking. Faint ambient street noise, possibly distant traffic or general outdoor sounds, is audible in the background;', 'moods': 'Informative, casual, engaging, and direct.'}]



VISTA at iteration 2: An 8-second video unfolds, starting outdoors on a lively, sunlit day. For the first 5.5 seconds, the camera captures a man in a static, chest-up shot. He is adorned with a full beard, a red baseball cap, black sunglasses featuring striking blue reflective lenses, a purple hoodie, and black over-ear headphones. Behind him, a gently blurred backdrop showcases a serene suburban street, complete with residential homes and abundant green foliage, all beneath a clear, light blue sky. The reflections visible in his sunglasses are flawlessly clean, ensuring no camera equipment or unwanted details are present to break the visual immersion. The man directly addresses the viewer, his voice clear and unaffected by any wind. A soft, consistent hum, suggesting distant street activity or general outdoor sounds, subtly enhances the natural ambiance without becoming a distraction. He asks, "Which comedian is known for their deadpan delivery?" Simultaneously, a sleek, black rectangular text box with white lettering gracefully appears from the bottom, positioned subtly at the screen's center-bottom, displaying "WHICH COMEDIAN IS KNOWN FOR THEIR DEADPAN DELIVERY?". After a short, expectant pause, he states, "Jeff Dye." The text overlay then smoothly transitions: the question fades out as "JEFF 'DYE'" fades in, maintaining its discreet size and placement. A subtle, knowing smile graces his lips as he concludes. At the 5.5-second mark, the video smoothly transitions from the live-action segment to the outro, utilizing a quick fade to black followed by a fade in. The final 2.5 seconds feature a pristine, minimalist white background. "Master of Puns" in black text, accompanied by a colorful 3D glasses emoji, is centered above a black rectangular "SUBSCRIBE!" button with a subtle grey border. A crisp 'pop' sound effect marks the video's conclusion.



VISTA at iteration 3: The 8-second video opens outdoors on a vibrant, sunlit day. For the initial 5.5 seconds, the camera presents a static, chest-up, eye-level view of a man. He sports a full beard, a red baseball cap, black sunglasses with striking blue reflective lenses, a purple hoodie, and black over-ear headphones. Behind him, a gently blurred scene of a serene suburban street, complete with residential homes and abundant green foliage, stretches beneath a clear, light blue sky. The reflections in his sunglasses are impeccably clean, free of any camera equipment or unwanted details, preserving visual immersion. The man addresses the viewer directly, his voice clear and unaffected by wind. A soft, consistent hum, hinting at distant street activity or general outdoor sounds, subtly enhances the natural ambiance without distraction. He asks, "Which comedian is known for their deadpan delivery?" Simultaneously, a sleek, black rectangular text box with white lettering gracefully emerges from the bottom, positioned subtly at the screen's center-bottom, displaying "WHICH COMEDIAN IS KNOWN FOR THEIR DEADPAN DELIVERY?". After a brief, expectant pause, he declares, "Jeff Dye." The text overlay then smoothly shifts: the question fades out as "JEFF 'DYE'" fades in, maintaining its discreet size and placement. A subtle, knowing smile plays on his lips as he finishes. At the 5.5-second mark, the video transitions from the live-action segment to the outro, marked by a quick fade to black followed by a fade in. The final 2.5 seconds showcase a pristine, minimalist white background. "Master of Puns" in black text, accompanied by a colorful 3D glasses emoji, is centered above a black rectangular "SUBSCRIBE!" button with a subtle grey border. A crisp 'pop' sound effect signals the video's conclusion.



VISTA at iteration 4: {'prompt_name': 'Trivia Master's Outdoor Riddle', 'prompt_description': 'The 8-second video opens outdoors on a vibrant, sunlit day. For the initial 5.5 seconds, the camera presents a static, chest-up, eye-level view of a man. The composition is artfully balanced, utilizing the shallow depth of field to enhance the subject's presence and create a visually engaging frame. He sports a full beard, a red baseball cap, black sunglasses with striking blue reflective lenses, a purple hoodie, and black over-ear headphones. Behind him, a gently blurred scene of a serene suburban street, complete with residential homes and abundant green foliage, stretches beneath a clear, light blue sky. The bright, sunny outdoor setting is bathed in a warm, inviting light, creating a subtly cinematic atmosphere. The reflections in his sunglasses are impeccably clean, free of any camera equipment or unwanted details, preserving visual immersion. A clear, warm, and resonant male voice speaks directly to the viewer. Faint, yet subtly varied, ambient street noise, including distant traffic and occasional natural outdoor sounds (e.g., a faint bird chirp or rustle of leaves), is audible in the background, adding to the immersive realism. The man asks, "Which comedian is known for their deadpan delivery?" Simultaneously, a sleek, black rectangular text box with white lettering gracefully emerges from the bottom, positioned subtly at the screen's center-bottom, displaying "WHICH COMEDIAN IS KNOWN FOR THEIR DEADPAN DELIVERY?". After a brief, expectant pause, he declares, "Jeff Dye." The text overlay then smoothly shifts: the question fades out as "JEFF 'DYE'" fades in, maintaining its discreet size and placement. A subtle, knowing smile plays on his lips as he finishes. At the 5.5-second mark, the video transitions from the live-action segment to the outro, marked by a quick fade to black followed by a fade in. The final 2.5 seconds showcase a pristine, minimalist white background. "Master of Puns" in black text, accompanied by a colorful 3D glasses emoji, is centered above a black rectangular "SUBSCRIBE!" button with a subtle grey border. A crisp 'pop' sound effect signals the video's conclusion.'}

Figure 13. Examples of prompts optimized by VISTA across iterations. Blue parts are updated.



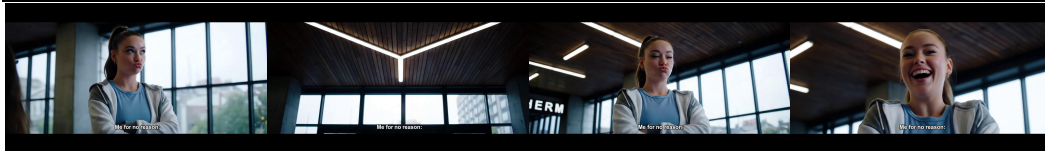
DP generates video with **gremlins moving backward fast without wooden rollercoaster**, which is physically non-sense.

VISTA fixes the issues with better visual fidelity: gremlins are moving forward and camera is backward.

Prompt: A rapid tracking shot of small, big-eared gremlins on a wooden rollercoaster in a midcentury theme park...



DP fails to cut between the **character's face shot and the interior ceiling shot**, nor repeat this transition.



VISTA fixes this with an even better camera focus and more natural character's expression.

Prompt: A short, humorous video depicting a young woman's mood rapidly shifting from bored and slightly annoyed to overtly joyful, triggered by a sudden and dramatic change in background music...

Figure 14. Comparison of Direct Prompting (DP) and VISTA with Veo 3. Top rows show Veo 3 failing to follow human instructions and generating unrealistic gremlin motions. Bottom rows highlight its inability to produce coherent multi-scene videos under DP. In contrast, VISTA corrects these issues with more realistic motions and smoother visual transitions.



Prompt: A spaceship entering hyperdrive, stars streaking past as it accelerates.

Self-Refine: The video is a strong visual representation of the prompt. The criticisms are minor and relate more to stylistic choices or common sci-fi tropes rather than a failure to meet the core request. It successfully conveys a spaceship accelerating into a hyperdrive state with stars streaking past.

VISTA (Visual Fidelity): The Normal Judge praises the high technical quality, clarity, sharp details, and consistent sci-fi aesthetic, noting the well-executed lighting and vibrant engine glow. They suggest minor improvements like more dynamic lighting and nuanced atmospheric distortion. The Negative Judge, however, critically points out **uninspired lighting** on the spaceship, and **static background galaxy, and overly uniform, artificial-looking streaking stars** lacking variation or parallax. They also mention the abruptness of transitions affecting fidelity. Upon review, the Negative Judge's assessment carries more weight. While the video is technically clean, the lighting on the ship does appear quite uniform, lacking the dynamic interplay of light and shadow expected in space...



Self-Refine: Overall, the generated video is **highly successful** in fulfilling the prompt. It accurately and visually compellingly depicts a spaceship entering hyperdrive with stars streaking past as it accelerates. There are no significant "wrongs" or misinterpretations of the prompt. The minor points above are more about stylistic choices or slight enhancements rather than fundamental errors.

VISTA (Motions and Dynamics): While the Normal Judge praises the smoothness of the ship's motion and dynamic star streaking, the Negative Judge correctly identifies a major directional flaw: **the spaceship moves vertically, which conflicts with viewer expectations of horizontal acceleration**. Additionally, the Negative Judge points out **the lack of micro-dynamics (e.g., rotational drift, buildup phases) and unrealistic exhaust behavior, which diminish the believability of motion**. These omissions outweigh the surface-level smoothness.



Prompt: The couple runs hand in hand to release a sky lantern, then watches it drift upward into the night sky, carried by the wind with the stars shining above.

Self-Refine: The video successfully portrays the actions of the couple running and releasing a lantern. However, it **fails significantly on the atmospheric detail of the "night sky with stars shining above."** The absence of stars is a major oversight that detracts from the prompt's intended mood and visual.

VISTA (Contextual Suitability): The scene struggles significantly with internal logical consistency, particularly regarding its temporal and environmental progression. An abrupt and jarring shift from a bright, late-afternoon setting to a deep, artificial blue night sky without any visual cues for time passing severely undermines immersion. This featureless 'night sky' further detracts from plausibility, failing to represent a natural environment. While individual actions are suitable, their placement within such an inconsistent context creates a sense of unreality.

VISTA (Visual Characters): While the two characters are central and necessary to the narrative, their portrayal lacks emotional depth and genuine interaction. Their expressions, particularly during the lantern release, appear somewhat static and posed, missing an opportunity to convey a deeper sense of wonder, intimacy, or shared experience. This underdevelopment makes them feel more like subjects performing actions rather than individuals experiencing a meaningful moment, limiting the scene's emotional resonance.

Figure 15. VISTA delivers high-level feedback on visual fidelity, contextual consistency, and emotional expression, surfacing nuanced flaws (e.g., static lighting, implausible scene transitions, and lack of character depth) that conventional LLM critiques often miss but are readily noticed by human judges.