
Video Models Can Reason with Verifiable Rewards

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Video diffusion models have made rapid progress in perceptual realism and tempo-
2 ral coherence, but they remain primarily optimized for plausible generation rather
3 than verifiable reasoning. This limitation is especially pronounced in tasks where
4 generated videos must satisfy explicit spatial, temporal, or logical constraints.
5 Inspired by the role of reinforcement learning with verifiable rewards (RLVR) in
6 reasoning-oriented language models, we introduce **VideoRLVR**, a practical recipe
7 for optimizing video diffusion models with rule-based feedback. VideoRLVR
8 formulates video reasoning as the generation of verifiable visual trajectories and
9 consists of an SDE-GRPO optimization backbone, dense decomposed rewards,
10 and an *Early-Step Focus* strategy for efficient training. The Early-Step Focus
11 strategy restricts policy optimization to the early denoising phase, reducing training
12 latency by about 40% while preserving performance. We evaluate VideoRLVR
13 on Maze, FlowFree, and Sokoban, three procedurally generated domains with
14 objective success criteria. Across these tasks, VideoRLVR consistently improves
15 over supervised fine-tuning baselines, with dense decomposed rewards proving
16 especially important in low-success-rate settings. Our RL-optimized model also
17 outperforms the evaluated proprietary and open-source video generation models
18 on these verifiable reasoning benchmarks. These results suggest that verifiable
19 RL can move video models beyond perceptual imitation toward more reliable
20 rule-consistent visual reasoning.

21 1 Introduction

22 Recent progress in large language models (LLMs) has reshaped the role of generative models from
23 content producers into increasingly capable reasoning systems [11, 32, 6]. A key intuition behind this
24 shift is that the model can externalize the problem-solving process by generating intermediate states
25 rather than only a final answer. This raises a natural question for video generation: *if language models*
26 *can reason through sequences of tokens, can video models reason through sequences of frames?*
27 Videos provide an appealing foundation for this idea, where each frame can represent an intermediate
28 visual state in a goal-directed process. In domains such as navigation [7], puzzle solving [14], and
29 embodied planning [27], a generated video can therefore be viewed not merely as motion synthesis,
30 but as a temporally ordered chain of visual states [39] that encodes a visual reasoning trajectory.

31 Despite this potential, current video diffusion models are still primarily optimized for perceptual
32 quality, temporal coherence, and plausible motion [13, 45, 36]. While large-scale video models
33 have begun to show signs of visual reasoning [39, 12, 37], these abilities remain difficult to elicit
34 reliably and verify under standard training objectives. The core challenge is the mismatch between
35 perceptual plausibility and objective correctness. Supervised fine-tuning (SFT) on ground-truth
36 solution videos can teach the model the visual form of valid trajectories, yet it does not directly
37 optimize the correctness of sampled outputs. As a result, models may imitate solution-like patterns
38 while failing to satisfy the underlying rules that make those solutions valid [9, 28]. This suggests an
39 analogy to reasoning-oriented LLMs where pre-training provides broad generative competence, SFT

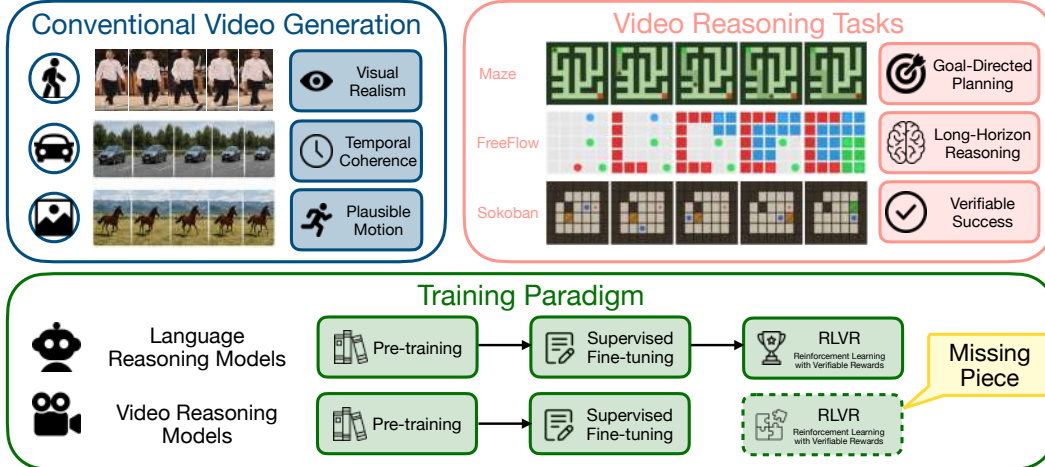


Figure 1: **Evolution towards verifiable video reasoning.** *Top:* Comparison between perception-focused generation and reasoning-intensive tasks. *Bottom:* We introduce **VideoRLVR**, the missing puzzle in the training paradigm for video reasoning models.

40 teaches the format of reasoning traces, Reinforcement Learning with Verifiable Rewards (RLVR) is
 41 the essential third stage required to optimize objective correctness.

42 In this work, we introduce **VideoRLVR**, a systematic recipe for applying reinforcement learning
 43 with verifiable rewards to video models. Our framework has four main components. First, we
 44 construct verifiable video reasoning data by generating solution trajectories with rule-based planners
 45 and aligning each logical transition with the video frame sequence. Second, we adopt an SDE-GRPO
 46 optimization backbone [24] for optimizing flow-matching video models. Third, we propose an
 47 *Early-Step Focus* strategy for efficient video RL. Instead of applying stochastic exploration and
 48 backpropagation across the entire denoising trajectory, this strategy concentrates optimization on the
 49 early denoising phase, where coarse structure and long-range planning are largely determined [38].
 50 Finally, we design dense decomposed rewards that break sparse task success into verifiable structural
 51 components, providing informative feedback even when full success is rare.

52 We evaluate our RLVR recipe on a multi-task suite designed for rule-based verification, including
 53 Maze, FlowFree, and Sokoban. Our experiments show that VideoRLVR improves video reasoning
 54 beyond supervised imitation. Across all three domains, the RL-optimized model consistently achieves
 55 higher success rates than the SFT checkpoint used to initialize training, with gains of 6.1%, 5.5%,
 56 and 3.2% on Maze, FlowFree, and Sokoban, respectively. Compared with continued supervised
 57 training, VideoRLVR yields larger gains on harder tasks, suggesting that verifiable rewards provide an
 58 optimization signal beyond what can be captured by imitation alone. We further evaluate VideoRLVR
 59 on the out-of-domain split of VBVR [37], where VideoRLVR shows improved transfer beyond
 60 the training domains. Our ablations further show that dense decomposed rewards are crucial in
 61 low-success-rate domains, and that Early-Step Focus reduces training time by about 40% while
 62 maintaining nearly the same performance. Finally, VideoRLVR outperforms several proprietary and
 63 open-source video generation models on our verifiable reasoning benchmarks, indicating that targeted
 64 verifiable RL can substantially improve the logical correctness of generated visual trajectories.

65 In summary, our contributions are as follows:

- 66 1. We introduce **VideoRLVR**, a reinforcement learning framework that optimizes video diffu-
 67 sion models with verifiable rewards, including dense decomposed reward functions to provide
 68 informative feedback for rule-verifiable visual trajectories.
- 69 2. We introduce a scalable training pipeline that combines rule-based trajectory generation, SDE-
 70 GRPO optimization, dense decomposed rewards, and an Early-Step Focus strategy that reduces
 71 training time by about 40% while preserving the performance.
- 72 3. We show that VideoRLVR improves over supervised fine-tuning and competitive proprietary and
 73 open-source video generation models on Maze, FlowFree, and Sokoban, while also demonstrating
 74 improved out-of-domain transfer on VBVR.

75 2 Related Work

76 **Reinforcement learning for diffusion and flow models.** Reinforcement learning has increasingly
77 been used to align diffusion and flow-based generative models with human preferences, perceptual
78 objectives, and task-specific rewards [42]. Prior work formulates denoising as a sequential decision
79 process and applies policy-gradient or preference-optimization methods to improve text-to-image
80 and video generation [2, 8, 35]. For flow-matching models, recent methods address the deterministic
81 nature of ODE sampling by introducing stochastic transitions or alternative preference objectives,
82 enabling likelihood-ratio or GRPO-style optimization [24, 43, 5]. Other extensions apply these
83 ideas to video or embodied objectives [1, 25]. However, existing work optimizes perceptual or
84 preference-based criteria such as aesthetics, text rendering, image fidelity, geometric consistency,
85 or motion quality [21, 22]. In contrast, our work studies reinforcement learning for *verifiable video*
86 *reasoning*, where rewards are computed from objective task rules and success depends on the logical
87 correctness of the generated visual trajectory.

88 **Reasoning in video generation models.** Recent work has begun to investigate whether video
89 generation models can serve as reasoning systems rather than only visual synthesizers. Large-scale
90 video models have shown emerging abilities on visual puzzles and sequential prediction tasks,
91 motivating the view that video generation can be interpreted as a chain of visual states or “chain of
92 frames” [39, 12, 18]. Benchmark efforts [37, 4, 44, 34] further evaluate video models on reasoning-
93 oriented tasks that require temporal consistency, spatial planning, or rule satisfaction. Other studies
94 analyze video models as world simulators or physical reasoners, highlighting both their potential
95 and their limitations in capturing causal and physical structure [3, 19, 27, 28, 47, 33]. These works
96 suggest that video models may contain useful visual reasoning priors, but also show that standard
97 generation objectives do not reliably produce rule-correct trajectories [12, 26]. Our work addresses
98 this gap by directly optimizing video models with verifiable rewards, using rule-based success criteria
99 rather than relying solely on supervised imitation or zero-shot generation.

100 **Verifiable reinforcement learning and reasoning models.** Reinforcement learning with verifiable
101 rewards has played an important role in recent progress on reasoning-oriented language models [11,
102 32, 6]. In these settings, the model is rewarded according to objective correctness signals, such as
103 mathematical equivalence, executable code tests, or rule-based verification, instead of only human
104 preference judgments [23, 46, 16, 17]. This paradigm is attractive because it provides scalable
105 supervision when outcomes can be automatically checked, which facilitates the development of
106 emerging behaviors like searching and backtracking [48, 41]. Our work extends this training from
107 language outputs to video trajectories. Whereas text reasoning is often verified by final-answer
108 correctness, video reasoning requires trajectory-level verification over visual, temporal, and process
109 constraints. We study how verifiable RL can optimize video diffusion models under these criteria.

110 3 Problem Formulation

111 **RLVR for Video Reasoning.** Following [39], we formulate video reasoning as a conditional
112 generation task where a model generates a temporal sequence of visual states whose transitions
113 and terminal state can be checked against task-specific rules. Given an initial image I_0 and a
114 textual instruction T , let $c = (I_0, T)$ denote the conditioning input. The model generates a video
115 $\mathbf{V} = \{I_0, I_1, \dots, I_{F-1}\}$, where F is the number of frames. Unlike standard video synthesis, which
116 primarily evaluates perceptual quality and temporal coherence, video reasoning requires the generated
117 sequence to satisfy task-specific correctness criteria. This formulation allows us to treat video
118 generation as a search for a valid visual trajectory conditioned on the initial state and instruction.

119 **Video Generation as a Markov Decision Process.** To apply reinforcement learning to flow-matching
120 video generation, we formulate the reverse denoising process as a Markov Decision Process (MDP)
121 over latent variables. This MDP is defined over denoising steps rather than reasoning steps, where the
122 reward is computed after the final video is decoded. At denoising step $k \in \{1, \dots, K\}$, the state is the
123 noisy video latent x_{t_k} at noise level t_k , and the action is the model velocity prediction $\hat{v}_\theta(x_{t_k}, t_k, c)$,
124 which determines the mean update of the next latent. Under the deterministic ODE solver, the
125 transition is given by $x_{t_{k+1}} = x_{t_k} + (t_{k+1} - t_k)a_k$. After the final denoising step, the decoded video
126 \mathbf{V} receives a verifier-derived reward $R(\mathbf{V}, c)$. A fundamental challenge in this formulation is that
127 standard flow matching employs a deterministic ODE solver, making it a deterministic function

128 of the initial noise x_1 . Under this deterministic solver, the next latent is a deterministic function
 129 of (x_{t_k}, c) , yielding no tractable stochastic transition density $\pi_\theta(x_{t_{k+1}} | x_{t_k}, c)$ for likelihood-ratio
 130 policy gradients. In Section 4, we address this by adopting an SDE-based formulation that introduces
 131 stochastic transitions compatible with flow-matching generation.

132 **Tasks.** To evaluate VideoRLVR across different reasoning domains, we instantiate our framework on
 133 three rule-verifiable visual reasoning domains: Maze, FlowFree, and Sokoban. We choose these tasks
 134 because they satisfy three properties: 1) solution correctness can be checked by rule-based verifiers,
 135 2) large-scale training and test instances can be generated, and 3) the tasks span different levels of
 136 reasoning complexity. Maze primarily tests spatial connectivity under explicit obstacle constraints,
 137 FlowFree requires globally consistent non-overlapping path connectivity and implicit constraints, and
 138 Sokoban introduces object interaction, irreversible transitions, and longer-horizon reasoning.

139 4 RLVR Recipe for Video Reasoning Models

140 We present **VideoRLVR**, a systematic recipe for optimizing video models with verifiable rewards.
 141 The recipe consists of four components: 1) rule-based data curation, 2) an SDE-GRPO optimization
 142 backbone, 3) dense decomposed rewards design, and 4) an Early-Step Focus optimization strategy.

143 4.1 Data Curation

144 Existing video reasoning datasets [44, 37] often lack the scale, task diversity, or fine-grained difficulty
 145 variation required to study RLVR for high-dimensional video generation. We therefore synthesize
 146 task instances with rule-based planners, producing conditioning inputs $c = (I_0, T)$ and valid video
 147 trajectories, together with metadata for automatic verification and reward computation. To align
 148 visual generation with symbolic reasoning, we map each discrete action to a unique frame transition
 149 $I_f \rightarrow I_{f+1}$. This action-to-frame mapping makes the generated video directly interpretable as a
 150 reasoning trajectory and enables the verifier to evaluate both terminal correctness and process-level
 151 validity. Dataset scale and task-specific generation details are provided in Section 5.

152 4.2 SDE-GRPO for Video Reasoning

153 GRPO [31] estimates relative advantages from groups of sampled outputs without training a separate
 154 critic, making it well suited for verifiable reward settings. However, standard flow-matching models
 155 generate samples with a deterministic ODE sampler, which does not provide a tractable stochastic
 156 transition density over denoising steps. Following Flow-GRPO [24], we convert the deterministic
 157 denoising dynamics into stochastic transitions with Gaussian log-probabilities.

158 **Stochastic denoising transitions.** For a discretized denoising schedule $\{t_k\}_{k=1}^K$, the SDE formula-
 159 tion defines a Gaussian transition:

$$\pi_\theta(x_{t_{k+1}} | x_{t_k}, c) = \mathcal{N}(x_{t_{k+1}}; \mu_\theta(x_{t_k}, t_k, c), \sigma_k^2 \mathbf{I}), \quad (1)$$

160 where $\mu_\theta(x_{t_k}, t_k, c)$ is the mean update induced by the model and σ_k^2 is the SDE transition variance.
 161 This stochastic transition enables closed-form log-probabilities and likelihood-ratio policy gradients.

162 **GRPO objective.** Given a group of G sampled videos for each condition, we compute verifier-
 163 derived rewards and normalize them within the group to obtain advantages A_i . For each sample i and
 164 denoising step k , we compute the dimension-normalized log-ratio:

$$\log \rho_{i,k} = \log \frac{\pi_\theta(x_{t_{k+1}}^{(i)} | x_{t_k}^{(i)}, c_i)}{\pi_{\text{old}}(x_{t_{k+1}}^{(i)} | x_{t_k}^{(i)}, c_i)} = -\frac{1}{2\sigma_k^2} \cdot \frac{1}{D} \sum_{d=1}^D \left[\left(x_{t_{k+1}}^{(i)} - \mu_\theta^{(i,k)} \right)_d^2 - \left(x_{t_{k+1}}^{(i)} - \mu_{\text{old}}^{(i,k)} \right)_d^2 \right], \quad (2)$$

165 where $\mu_\theta^{(i,k)} = \mu_\theta(x_{t_k}^{(i)}, t_k, c_i)$, $\mu_{\text{old}}^{(i,k)} = \mu_{\text{old}}(x_{t_k}^{(i)}, t_k, c_i)$, and D is the number of latent elements.
 166 The policy loss uses PPO-style clipping:

$$\mathcal{L}_{\text{policy}} = -\mathbb{E}_{i,k} [\min(\rho_{i,k} A_i, \text{clip}(\rho_{i,k}, 1 - \varepsilon, 1 + \varepsilon) A_i)]. \quad (3)$$

167 We additionally regularize the policy against the reference model with a closed-form KL penalty:

$$\mathcal{L}_{\text{KL}} = \mathbb{E}_k \left[\frac{1}{D} \frac{\|\mu_\theta - \mu_{\text{ref}}\|_2^2}{2\sigma_k^2} \right]. \quad (4)$$

168 The final objective is $\mathcal{L}_{\text{VideoRLVR}} = \mathcal{L}_{\text{policy}} + \beta \mathcal{L}_{\text{KL}}$, where β controls the strength of regularization.

169 **4.3 Early-Step Focus for Efficient Video RL**

170 Video RL is substantially more expensive than text RL because each rollout requires generating and
 171 backpropagating through high-dimensional spatio-temporal latents. A full SDE-GRPO update over
 172 all K denoising steps therefore incurs large memory and time costs. However, not all denoising steps
 173 contribute equally to the reasoning objective. Early high-noise steps are primarily responsible for
 174 coarse layout, object placement, and long-range structure, whereas later low-noise steps mainly refine
 175 local appearance and consolidate the generation into a specific visual trajectory [38].

176 Motivated by this observation, we introduce *Early-Step Focus*. During RL optimization, we sample
 177 the full denoising trajectory for generation and reward evaluation, but restrict stochastic perturbation,
 178 log-probability computation, and gradient backpropagation to the first $L < K$ denoising steps.
 179 This creates an efficient exploration-exploitation trade-off: early denoising steps receive stochastic
 180 perturbations and policy-gradient updates for high-level reasoning, while later steps preserve the
 181 generative prior and refine visual details. The policy loss becomes:

$$\mathcal{L}_{\text{ESF}} = -\mathbb{E}_{i,k \leq L} [\min(\rho_{i,k} A_i, \text{clip}(\rho_{i,k}, 1 - \varepsilon, 1 + \varepsilon) A_i)] + \beta \mathcal{L}_{\text{KL}}^{k \leq L}. \quad (5)$$

182 In our experiments, we use $K = 20$ denoising steps and $L = 10$ early steps. This reduces training
 183 latency by about 40% while preserving reasoning performance, suggesting that the early denoising
 184 phase carries most of the reward-relevant structural signal.

185 **4.4 Reward Design for Video Reasoning**

186 Verifiable reasoning tasks provide objective rules that can be converted into dense reward signals.
 187 Instead of using only a binary success indicator, we decompose each task into structural components
 188 that measure partial progress toward a valid solution. This is especially important in low-success-rate
 189 domains, where most sampled videos fail completely and binary rewards provide little variation
 190 within a GRPO group.

191 **Task-aware Reward Function.** We use a task-aware reward function for joint training across
 192 heterogeneous domains. For each conditioning input c , the dispatcher identifies the task $\mathcal{T}(c) \in$
 193 $\{\text{Maze}, \text{FlowFree}, \text{Sokoban}\}$ and evaluates the generated video with the corresponding reward:

$$R(\mathbf{V}, c) = R_{\mathcal{T}(c)}(\mathbf{V}, c). \quad (6)$$

194 This allows mixed-task RL batches while preserving task-specific verification criteria.

195 **Dense Reward Formulations.** For each task, we decompose the global objective into measurable
 196 rule-based components:

- 197 • **Maze.** We define the reward as: $R_{\text{maze}} = c_{\text{conn}} \cdot w_{\text{wall}}$, where c_{conn} measures start-to-goal path
 198 connectivity and w_{wall} penalizes wall violations. The multiplicative form couples connectivity
 199 with wall consistency, discouraging connected paths that violate maze constraints.
- 200 • **FlowFree.** We combine four structural metrics: $R_{\text{ff}} = \lambda_{\text{valid}} p_{\text{valid}} + \lambda_{\text{ep}} \text{ep}_{\text{pres}} + \lambda_{\text{conn}} c_{\text{conn}} +$
 201 $\lambda_{\text{fill}} \text{fill}_{\text{rate}}$, where p_{valid} measures endpoint-to-endpoint path validity, ep_{pres} measures preservation
 202 of the given endpoints, c_{conn} measures 4-connected color regions, and $\text{fill}_{\text{rate}}$ measures grid
 203 coverage by valid path colors. The weights $\lambda_{\text{valid}}, \lambda_{\text{ep}}, \lambda_{\text{conn}}, \lambda_{\text{fill}}$ balance the relative importance
 204 of these components. In our experiments, we set them to 0.15, 0.35, 0.30, and 0.20, respectively.
- 205 • **Sokoban.** We use a combination of final-state and process-validity rewards: $R_{\text{sok}} = \lambda_{\text{state}} R_{\text{state}} +$
 206 $\lambda_{\text{proc}} R_{\text{proc}}$, where R_{state} measures box placement on target cells and R_{proc} measures the fraction of
 207 valid transitions under Sokoban movement rules. The weights λ_{state} and λ_{proc} balance final-state
 208 correctness and process validity. We use $\lambda_{\text{state}} = \lambda_{\text{proc}} = 0.5$ in all experiments.

209 **5 Experiments**

210 In this section, we evaluate **VideoRLVR** from two perspectives. First, we compare against supervised
 211 fine-tuning and competitive video generation baselines on three rule-verifiable reasoning domains:
 212 Maze, FlowFree, and Sokoban. Then, we test transfer beyond the training domains using the out-of-
 213 domain split of VBVR [37]. Together, these experiments assess whether verifiable RL improves both
 214 in-domain rule-based correctness and out-of-domain visual reasoning behavior.

Table 1: Comparison of our method with other state-of-the-art method. We report both Success Rate (SR) and GT alignment metrics: Precision (Prec), Recall (Rec), and F1. **Bold** indicates the best and underlined indicates second best.

| Model | Maze | | | | FlowFree | | | | Sokoban | | | |
|---------------------------|------|------|------|-------------|----------|------|------|------------|---------|------|------|------------|
| | Prec | Rec | F1 | SR | Prec | Rec | F1 | SR | Prec | Rec | F1 | SR |
| <i>Proprietary Models</i> | | | | | | | | | | | | |
| Sora 2 | 15.8 | 17.2 | 16.5 | 3.1 | 10.8 | 5.1 | 5.8 | 0.0 | 8.5 | 4.8 | 5.4 | 0.0 |
| Kling V3 | 24.8 | 15.7 | 19.2 | 23.5 | 18.8 | 2.7 | 4.7 | 0.0 | 5.7 | 2.7 | 3.7 | 0.0 |
| Veo 3.1 | 22.8 | 18.1 | 20.2 | 26.0 | 23.9 | 4.7 | 7.5 | <u>4.0</u> | 22.2 | 6.0 | 9.4 | 0.0 |
| <i>Open-Source Models</i> | | | | | | | | | | | | |
| CogVideoX1.5 | 13.3 | 10.8 | 11.9 | 0.0 | 18.7 | 2.2 | 3.9 | 0.0 | 3.2 | 0.3 | 0.5 | 0.0 |
| HunyuanVideo | 17.3 | 11.4 | 13.8 | 2.2 | 12.5 | 2.9 | 4.8 | 0.0 | 8.2 | 2.7 | 3.2 | 0.0 |
| Wan2.2-TI2V-5B | 18.3 | 12.2 | 14.6 | 0.0 | 17.4 | 2.0 | 3.4 | 0.0 | 4.1 | 0.7 | 1.0 | 0.0 |
| <i>SFT Models</i> | | | | | | | | | | | | |
| Wan-R1 | 20.9 | 65.6 | 31.7 | 31.9 | 20.9 | 3.6 | 6.1 | 0.0 | 7.7 | 2.1 | 3.3 | 0.0 |
| VBVR-Wan2.2 | 62.7 | 77.8 | 69.4 | 60.8 | 17.9 | 5.6 | 8.5 | 1.7 | 16.2 | 1.7 | 3.1 | 0.0 |
| SFT Epoch 5 | 80.2 | 83.0 | 81.6 | 66.1 | 42.8 | 42.2 | 42.4 | 2.4 | 33.6 | 11.9 | 17.6 | 2.9 |
| SFT Epoch 10 | 80.4 | 85.1 | 82.7 | <u>69.0</u> | 43.1 | 42.5 | 42.8 | 2.5 | 32.8 | 11.6 | 17.1 | 2.7 |
| <i>RL Model</i> | | | | | | | | | | | | |
| Ours | 82.1 | 86.9 | 84.4 | 72.2 | 44.3 | 43.8 | 44.0 | 7.9 | 34.0 | 12.5 | 29.4 | 6.1 |

215 5.1 Experimental Setup

216 **Dataset.** We train and evaluate on a multi-task suite of three procedurally generated reasoning
217 domains: Maze, FlowFree, and Sokoban. To prevent the model from overfitting to specific visual
218 features, we apply varied color themes across the dataset, encouraging the model to rely on structural
219 invariants. Each sample consists of an input image, a task instruction, and an 81-frame ground-truth
220 video at 480×832 resolution. The total training dataset consists of 30,000 samples (1,0000 per task).
221 For the test set, we maintain a held-out set of 3,000 samples (1,000 per task) generated with disjoint
222 random seeds. Dataset construction details are provided in Section A.1.

223 **Base Model and SFT Baseline.** We use Wan2.2-TI2V-5B [36], a state-of-the-art video generation
224 model, as our base model. It generates $F = 81$ frames at 480×832 resolution. We first establish an
225 SFT baseline by training the model on ground-truth solution videos using the standard flow matching
226 objective. This SFT checkpoint provides the necessary perceptual and structural prior for the model,
227 serving as both the initial policy and the reference policy π_{ref} for RL optimization.

228 **Baselines.** To evaluate the effectiveness of our method, we compare our model with competitive
229 proprietary and open-source video generation baselines. For proprietary models, we use Sora 2 [29],
230 Kling V3 [20], and Veo 3.1 [39]. For open-source models, we compare with Wan2.2-TI2V-5B [36],
231 CogVideoX1.5-5B-I2V [13], and HunyuanVideo-I2V [40]. We also compare with specialized SFT-
232 based video reasoning models, including Wan-R1 [44] and VBVR-Wan2.2 [37]. Wan-R1 adopts
233 the same base model as ours and is trained on the Maze and Sokoban domains with LoRA [15].
234 VBVR-Wan2.2 utilizes Wan2.2-I2V-A14B [36] and is trained on the VBVR dataset with LoRA.

235 **Training Configuration.** We train the SFT baseline for 5 epochs with a learning rate of 1×10^{-5} .
236 For VideoRLVR, we initialize from the SFT checkpoint and train on the same training set for 1 epoch
237 using SDE-GRPO as the optimization backbone. We use group size $G = 16$, $T = 20$ denoising
238 steps, learning rate 5×10^{-6} , and KL coefficient $\beta = 0.04$. Following the Early-Step Focus strategy
239 in Section 4.3, backpropagation and SDE injection is restricted to the first $L = 10$ steps of the
240 denoising trajectory. All training experiments are conducted on 8 NVIDIA B200 GPUs.

241 **Evaluation.** We evaluate the results using two complementary metric families: 1) trajectory
242 alignment metrics, including Precision (Prec), Recall (Rec), and F1, which measure pixel-, cell-,
243 or action-level alignment with the reference solution, and 2) symbolic success rate, which verifies
244 whether the video satisfies the underlying task rules. Evaluation details are listed in Section A.2

245 5.2 Main Results

246 Table 1 compares VideoRLVR with supervised baselines and competitive video generation models
247 on our verifiable reasoning benchmarks.

248 **RLVR consistently outperforms supervised baselines.** VideoRLVR yields consistent improvements
 249 across all three reasoning domains. Compared with the SFT Epoch 5 checkpoint used to initialize
 250 RL training, VideoRLVR improves success rate by 6.1% on Maze, 5.5% on FlowFree, and 3.2% on
 251 Sokoban. Notably, VideoRLVR also significantly surpasses the performance of recent state-of-the-art
 252 closed-source models on visual reasoning tasks, validating the efficacy of verifiable reinforcement
 253 learning in domains where generic video pre-training remains insufficient for complex logical tasks.

254 **Superior scaling on high-complexity tasks.** To isolate the specific advantages of RLVR over
 255 extended supervised learning, we evaluate a stronger SFT baseline (SFT Epoch 10), repre-
 256 senting the result of conducting further supervised training on the same Epoch 5 checkpoint.
 257 As shown in Table 1 and Figure 2, VideoRLVR
 258 is more robust than continued SFT as task diffi-
 259 culty increases. Within the Maze domain, RLVR
 260 establishes a 3.2% margin over the SFT Epoch
 261 10 checkpoint and shows less degradation when
 262 the scale of the maze increases. On FlowFree,
 263 VideoRLVR improves over SFT Epoch 10 by
 264 5.4%, while continued SFT provides little im-
 265 provement over the Epoch 5 checkpoint. On
 266 Sokoban, continued SFT slightly degrades per-
 267 formance, whereas VideoRLVR improves over
 268 SFT Epoch 10 by 3.4%. These trends suggest
 269 that verifiable rewards provide an optimization
 270 signal that is not captured by additional imita-
 271 tion training alone.

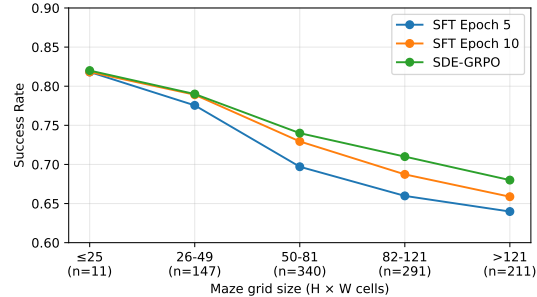


Figure 2: Success rate of different grid size for maze. n represents the number of samples falling into this range

272 5.3 Comparison with LLMs

273 To determine if our reasoning domains can be
 274 solved by language reasoning alone, we bench-
 275 mark frontier LLMs on the Maze task. Table 2
 276 presents the results of state-of-the-art models, in-
 277 cluding GPT-5.5 Pro [30] and Gemini 3.1 Pro [10],
 278 compared against our RLVR-optimized video
 279 model. Despite their sophisticated reasoning capa-
 280 bilities in textual domains, LLMs exhibit a sharp
 281 performance decay in maze tasks. This divergence
 282 highlights a representation bottleneck: while LLMs must reason over a tokenized rendering of the
 283 maze, our video model operates directly on a visual latent space that inherently preserves the visual
 284 topological relationships necessary for complex visual reasoning. These results suggest that, for
 285 visual reasoning, directly generating and optimizing visual trajectories can be more effective than
 286 solving the task through language-token representations alone.

Table 2: Comparison with LLMs on maze tasks.

| Model | Maze | | | |
|------------------|-------------|-------------|-------------|-------------|
| | Prec | Rec | F1 | SR |
| GPT 4o | 11.7 | 13.0 | 12.3 | 0.0 |
| GPT 5.5 Pro | 76.0 | 70.1 | 72.9 | 66.0 |
| Gemini 2.5 Flash | 11.2 | 10.5 | 10.9 | 0.0 |
| Gemini 3.1 Pro | 26.8 | 27.0 | 26.9 | 23.0 |
| Ours | 82.1 | 86.9 | 84.4 | 72.2 |

287 5.4 OOD Results

288 To evaluate whether Video-
 289 RLVR transfers beyond
 290 the training domains, we
 291 test our model on the out-of-
 292 domain split of VBVR [37].
 293 This benchmark covers mul-
 294 tiple reasoning categories
 295 and therefore provides a
 296 broader test of whether
 297 VideoRLVR improves gen-
 298 eral video reasoning behavior beyond Maze, FlowFree, and Sokoban. As shown in Table 3, VideoR-
 299 LVR substantially improves over the 5B baseline, increasing the average score from 26.2 to 60.2 with
 300 gains across all VBVR-ODD categories. VideoRLVR also performs competitively with the larger 14B
 301 VBVR-Wan2.2 model, achieving a similar average score despite using a smaller 5B backbone and

Table 3: **OOD evaluation on VBVR.** We report average performance and category-wise scores on the VBVR-ODD split.

| Model | Avg. | Abst. | Know. | Perc. | Spat. | Trans. |
|-------------------|-------------|-------------|-------------|-------------|-------------|-------------|
| <i>5B Models</i> | | | | | | |
| CogVideoX1.5 | 26.2 | 28.1 | 23.5 | 25.0 | 25.4 | 28.2 |
| VideoRLVR | 60.2 | 65.5 | 62.0 | 59.7 | 58.8 | 58.2 |
| <i>14B Models</i> | | | | | | |
| Wan2.2-I2V-A14B | 32.9 | 40.5 | 30.8 | 34.3 | 23.6 | 30.7 |
| VBVR-Wan2.2 | 61.0 | 76.8 | <u>57.2</u> | <u>54.7</u> | 61.8 | 61.5 |

302 much less training data. These results suggest that VideoRLVR learns transferable visual reasoning
 303 ability that generalizes beyond the generated training tasks.

304 6 Analysis

305 In this section, we analyze the main components of VideoRLVR, including GRPO group size,
 306 Early-Step Focus, KL regularization, dense reward design, and qualitative generation behavior.

307 6.1 Ablation Study

308 **Impact of Group Size.** In GRPO, the group size G affects the stability of group-relative ad-
 309 vantage estimation. We investigate the impact of this hyperparameter within the Maze reasoning
 310 domain, as shown in Figure 3. Our results indi-
 311 cate that performance scales positively with
 312 group size, primarily due to the stabilization of
 313 the reward distribution’s statistics. While the
 314 expected sample standard deviation of rewards
 315 is a property of the policy’s diversity, a small
 316 group size (e.g., $G \leq 4$) provides a noisy and
 317 often biased estimate of this value. This leads
 318 to significant fluctuations in the advantage calcu-
 319 lation, as the group mean fails to accurately
 320 represent the current policy’s performance level.
 321 Increasing the group size to $G = 16$ provides a
 322 more stable comparison set for estimating relative advantages, which improves training stability in
 323 our experiments. However, we observe diminishing returns beyond this point. Furthermore, because
 324 video generation remains a significant computational bottleneck, scaling G entails a linear increase in
 325 rollout time and VRAM overhead. We therefore use $G = 16$ as a practical trade-off between
 326 advantage-estimation stability and computational cost.

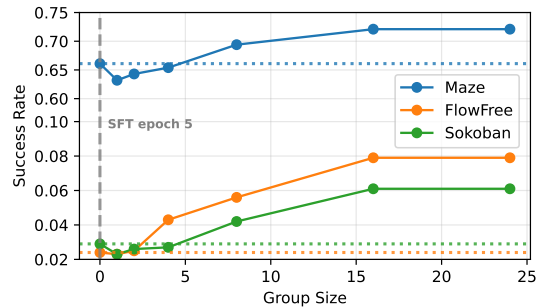


Figure 3: Scaling results with group size.

327 **Early-Step Focus.** To validate the efficacy of our Early-Step Focus strategy, we conduct a controlled
 328 experiment within the Maze domain. We fix the total inference budget at $T = 20$ denoising steps
 329 and compare the reasoning performance when the gradient and noise injection is calculated
 330 over the full trajectory ($L = 20$) versus the
 331 first $L = 10$ steps. As illustrated in Table 4,
 332 the success rates and F1 scores remain nearly
 333 unchanged, while training time is substantially
 334 reduced. This suggests that the early denoising steps carry much of the reward-relevant structural
 335 signal for visual reasoning. Because the later denoising steps primarily govern local textural re-
 336 finement, they contribute less to the verifier-derived reasoning objective in this setting. Restricting
 337 backpropagation and noise injection to these early steps thus serves as a computationally efficient
 338 optimization path, significantly reducing the training time without degrading the performance.

Table 4: Comparisons of computing over full denoising steps and early step focus.

| Gradient Steps | F1 | SR | Time / step |
|---------------------|------|------|-------------|
| 20 (Full) | 84.6 | 72.3 | 156 s |
| 10 (Early 10 steps) | 84.4 | 72.2 | 93.5 s |

340 **KL Constraint.** The KL-divergence constraint is essential for maintaining the model’s generative
 341 prior. We show a qualitative example in Section B.1. As shown in Figure 5, removing the KL penalty
 342 ($\beta = 0$) at an early stage of GRPO optimization can lead to reward-hacking behavior. Without this
 343 regularization, the model the model may produce visually implausible patterns that satisfy parts
 344 of the verifier while degrading generation quality. Implementing a constant penalty of $\beta = 0.04$
 345 successfully anchors the optimization to the original quality, ensuring that improvements in logical
 346 success are achieved without sacrificing the model’s inherent visual plausibility.

347 6.2 Reward Design

348 To evaluate the necessity of our dense de-
 349 composed reward design, we investigate
 350 the efficacy of a sparse reward ($R \in \{0, 1\}$)
 351 based exclusively on success rate. This ab-

Table 5: Training with a sparse binary success reward.

| Steps | 0 | 200 | 400 | 600 | 800 | 1000 |
|----------|------|------|------|------|------|------|
| Maze | 66.1 | 67.2 | 68.9 | 70.1 | 71.5 | 72.9 |
| FlowFree | 2.4 | 2.3 | 2.5 | 2.4 | 2.6 | 2.5 |

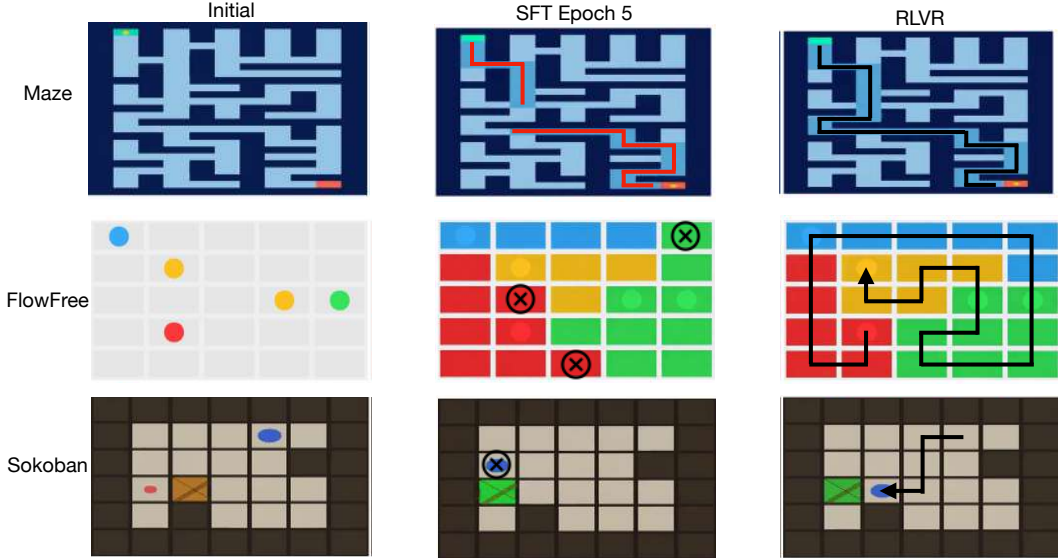


Figure 4: Qualitative case study across three reasoning domains. We compare generations from the SFT baseline (Epoch 5) and the final RLVR-optimized model.

352 lation aims to determine if binary feedback
 353 is sufficient across varying levels of task complexity.

354 Our results, shown in Table 5, reveal that sparse success rewards behave differently across domains.
 355 In domains like maze, where the baseline model already achieves a decent success rate, the sparse
 356 reward signal is sufficient to provide an informative gradient. The model is able to encounter success
 357 frequently enough during group rollouts to differentiate between advantageous and disadvantageous
 358 trajectories. Conversely, on high-complexity tasks like FlowFree and Sokoban, where the initial
 359 success rate is near-zero, the sparse reward provides little useful signal. In these environments,
 360 the model suffers from extreme gradient sparsity. Since success is rarely encountered within the
 361 group rollout G , the advantage estimates remain uninformative. This underscores a critical cold-start
 362 problem in video RL: while binary success is the ultimate goal, it is an insufficient signal for exploring
 363 high-dimensional latent space from a low-performance starting point. Dense decomposed rewards
 364 address this issue by providing partial credit for intermediate structural properties, giving the policy
 365 useful feedback before full success becomes frequent.

366 6.3 Case Study

367 Figure 4 provides qualitative examples across Maze, FlowFree, and Sokoban. The SFT baseline often
 368 captures the visual format of each domain, such as drawing paths, coloring grids, or rendering objects,
 369 but it can fail to satisfy the task rules. For example, SFT outputs may contain disconnected paths,
 370 inconsistent color connectivity, or invalid object transitions. In the Sokoban example, the SFT model
 371 produces a visually plausible but invalid shortcut rather than a valid sequence of box-pushing actions.
 372 In contrast, the VideoRLVR-optimized model more consistently satisfies the symbolic constraints
 373 checked by our verifiers. Across the shown examples, it produces connected paths, more coherent
 374 grid solutions, and more valid object transitions while preserving the overall visual structure of the
 375 task. These qualitative results support the quantitative findings: verifiable RL improves rule-based
 376 correctness beyond what is achieved by supervised imitation alone.

377 7 Conclusion

378 This work studies reinforcement learning with verifiable rewards for video reasoning and introduces
 379 **VideoRLVR**, a practical recipe for optimizing video reasoning models. By combining rule-verifiable
 380 data generation, an SDE-GRPO optimization backbone, dense decomposed rewards, and Early-Step
 381 Focus, VideoRLVR addresses the gap between perceptual video synthesis and task-level logical

382 correctness. Our experiments show that supervised fine-tuning provides an important visual and
383 structural prior, but can plateau or degrade on harder reasoning tasks when optimized only through
384 imitation. In contrast, verifiable RL improves success rates across Maze, FlowFree, and Sokoban,
385 with dense decomposed rewards proving especially useful in low-success-rate domains. We further
386 show that Early-Step Focus reduces training time by about 40% with little observed loss in reasoning
387 performance. Overall, our results suggest that verifiable RL can substantially improve the logical
388 correctness of generated videos, enabling open-source video models to outperform stronger general-
389 purpose video generation baselines on visual reasoning benchmarks.

390 **Limitations**

391 **Broader Impacts**

392 References

- 393 [1] Zhaochong An, Orest Kupyn, Théo Uscidda, Andrea Colaco, Karan Ahuja, Serge Belongie,
394 Mar Gonzalez-Franco, and Marta Tintore Gazulla. Vggrpo: Towards world-consistent video
395 generation with 4d latent reward. *arXiv preprint arXiv:2603.26599*, 2026.
- 396 [2] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion
397 models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.
- 398 [3] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Leo Jing, David Schnurr,
399 Joe Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators.
400 *OpenAI Blog*, 1(8):1, 2024.
- 401 [4] Zefan Cai, Haoyi Qiu, Tianyi Ma, Haozhe Zhao, Gengze Zhou, Kung-Hsiang Huang, Parisa
402 Kordjamshidi, Minjia Zhang, Wen Xiao, Jiuxiang Gu, et al. Mmgr: Multi-modal generative
403 reasoning. *arXiv preprint arXiv:2512.14691*, 2025.
- 404 [5] Wentse Chen, Shiyu Huang, Yuan Chiang, Tim Pearce, Wei-Wei Tu, Ting Chen, and Jun Zhu.
405 Dgpo: discovering multiple strategies with diversity-guided policy optimization. In *Proceedings*
406 *of the AAAI conference on artificial intelligence*, volume 38, pages 11390–11398, 2024.
- 407 [6] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit
408 Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the
409 frontier with advanced reasoning, multimodality, long context, and next generation agentic
410 capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- 411 [7] Yifei Dong, Fengyi Wu, Yilong Dai, Lingdong Kong, Guangyu Chen, Xu Zhu, Qiyu Hu, Tianyu
412 Wang, Johnalbert Garnica, Feng Liu, et al. Language-conditioned world modeling for visual
413 navigation. *arXiv preprint arXiv:2603.26741*, 2026.
- 414 [8] Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel,
415 Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpok: Reinforcement learning
416 for fine-tuning text-to-image diffusion models. *Advances in Neural Information Processing*
417 *Systems*, 36:79858–79885, 2023.
- 418 [9] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel,
419 Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature*
420 *Machine Intelligence*, 2(11):665–673, 2020.
- 421 [10] Google DeepMind. Gemini 3.1 Pro Model Card. [https://deepmind.google/models/
422 model-cards/gemini-3-1-pro/](https://deepmind.google/models/model-cards/gemini-3-1-pro/), February 2026.
- 423 [11] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu,
424 Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in
425 llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 426 [12] Ziyu Guo, Xinyan Chen, Renrui Zhang, Ruichuan An, Yu Qi, Dongzhi Jiang, Xiangtai Li,
427 Manyuan Zhang, Hongsheng Li, and Pheng-Ann Heng. Are video models ready as zero-shot
428 reasoners? an empirical study with the mme-cof benchmark. *arXiv preprint arXiv:2510.26802*,
429 2025.
- 430 [13] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale
431 pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*,
432 2022.
- 433 [14] Sepidehsadat Sepid Hossieni, Mohammad Amin Shabani, Saghar Irandoust, and Yasutaka
434 Furukawa. Puzzlefusion: Unleashing the power of diffusion models for spatial puzzle solving.
435 *Advances in Neural Information Processing Systems*, 36:9574–9597, 2023.
- 436 [15] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Liang
437 Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *Iclr*, 1(2):3,
438 2022.

- 439 [16] Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xiangyu Zhang, and Heung-Yeung Shum.
440 Open-reasoner-zero: An open source approach to scaling up reinforcement learning on the base
441 model. *arXiv preprint arXiv:2503.24290*, 2025.
- 442 [17] Yixu Huang, Tinghui Zhu, and Muhao Chen. Learning adaptive reasoning paths for efficient
443 visual reasoning. *arXiv preprint arXiv:2604.14568*, 2026.
- 444 [18] Ziqi Huang, Ning Yu, Gordon Chen, Haonan Qiu, Paul Debevec, and Ziwei Liu. Vchain:
445 Chain-of-visual-thought for reasoning in video generation. *arXiv preprint arXiv:2510.05094*,
446 2025.
- 447 [19] Bingyi Kang, Yang Yue, Rui Lu, Zhijie Lin, Yang Zhao, Kaixin Wang, Gao Huang, and Jiashi
448 Feng. How far is video generation from world model: A physical law perspective. *arXiv*
449 *preprint arXiv:2411.02385*, 2024.
- 450 [20] Kling AI. All you need to know about kling video 3.0. [https://kling.ai/blog/
451 kling-video-3-0-ai-director-features-guide](https://kling.ai/blog/kling-video-3-0-ai-director-features-guide), February 2026.
- 452 [21] Rui Li, Yuanzhi Liang, Ziqi Ni, Haibing Huang, Chi Zhang, and Xuelong Li. Growing with the
453 generator: Self-paced grpo for video generation. *arXiv preprint arXiv:2511.19356*, 2025.
- 454 [22] Yuming Li, Yikai Wang, Yuying Zhu, Zhongyu Zhao, Ming Lu, Qi She, and Shanghang Zhang.
455 Branchgrpo: Stable and efficient grpo with structured branching in diffusion models. *arXiv*
456 *preprint arXiv:2509.06040*, 2025.
- 457 [23] Zhong-Zhi Li, Duzhen Zhang, Ming-Liang Zhang, Jiaxin Zhang, Zengyan Liu, Yuxuan Yao,
458 Haotian Xu, Junhao Zheng, Pei-Jie Wang, Xiuyi Chen, et al. From system 1 to system 2: A
459 survey of reasoning large language models. *arXiv preprint arXiv:2502.17419*, 2025.
- 460 [24] Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan,
461 Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. *arXiv*
462 *preprint arXiv:2505.05470*, 2025.
- 463 [25] Xiao Liu, Yifan Zhou, Fabian Weigend, Shubham Sonawani, Shuhei Ikemoto, and Heni Ben
464 Amor. Diff-control: A stateful diffusion-based policy for imitation learning. In *2024 IEEE/RSJ*
465 *International Conference on Intelligent Robots and Systems (IROS)*, pages 7453–7460. IEEE,
466 2024.
- 467 [26] Yang Luo, Xuanlei Zhao, Baijiong Lin, Lingting Zhu, Liyao Tang, Yuqi Liu, Ying-Cong Chen,
468 Shengju Qian, Xin Wang, and Yang You. V-reasonbench: Toward unified reasoning benchmark
469 suite for video generation models. *arXiv preprint arXiv:2511.16668*, 2025.
- 470 [27] Zhiting Mei, Tenny Yin, Ola Shorinwa, Apurva Badithela, Zhonghe Zheng, Joseph Bruno, Madis-
471 on Bland, Lihan Zha, Asher Hancock, Jaime Fernández Fisac, et al. Video generation models in
472 robotics-applications, research challenges, future directions. *arXiv preprint arXiv:2601.07823*,
473 2026.
- 474 [28] Saman Motamed, Laura Culp, Kevin Swersky, Priyank Jaini, and Robert Geirhos. Do gener-
475 ative video models understand physical principles? In *Proceedings of the IEEE/CVF Winter*
476 *Conference on Applications of Computer Vision*, pages 948–958, 2026.
- 477 [29] OpenAI. Sora 2 system card. <https://openai.com/index/sora-2-system-card/>,
478 September 2025.
- 479 [30] OpenAI. GPT-5.5 System Card. <https://openai.com/index/gpt-5-5-system-card/>,
480 April 2026.
- 481 [31] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
482 Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical
483 reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- 484 [32] Aaditya Singh, Adam Fry, Adam Perelman, Adam Tart, Adi Ganesh, Ahmed El-Kishky, Aidan
485 McLaughlin, Aiden Low, AJ Ostrow, Akhila Ananthram, et al. Openai gpt-5 system card. *arXiv*
486 *preprint arXiv:2601.03267*, 2025.

- 487 [33] Selena Song, Ziming Xu, Zijun Zhang, Kun Zhou, Jiaxian Guo, Lianhui Qin, and Biwei
488 Huang. Learning plug-and-play memory for guiding video diffusion models. *arXiv preprint*
489 *arXiv:2511.19229*, 2025.
- 490 [34] Jingqi Tong, Yurong Mou, Hangcheng Li, Mingzhe Li, Yongzhuo Yang, Ming Zhang, Qiguang
491 Chen, Tianyi Liang, Xiaomeng Hu, Yining Zheng, et al. Thinking with video: Video generation
492 as a promising multimodal reasoning paradigm. *arXiv preprint arXiv:2511.04570*, 2025.
- 493 [35] Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam,
494 Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using
495 direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer*
496 *Vision and Pattern Recognition*, pages 8228–8238, 2024.
- 497 [36] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwei Yu,
498 Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative
499 models. *arXiv preprint arXiv:2503.20314*, 2025.
- 500 [37] Maijunxian Wang, Ruisi Wang, Juyi Lin, Ran Ji, Thaddäus Wiedemer, Qingying Gao, Dezhi
501 Luo, Yaoyao Qian, Lianyu Huang, Zelong Hong, et al. A very big video reasoning suite. *arXiv*
502 *preprint arXiv:2602.20159*, 2026.
- 503 [38] Ruisi Wang, Zhongang Cai, Fanyi Pu, Junxiang Xu, Wanqi Yin, Maijunxian Wang, Ran
504 Ji, Chenyang Gu, Bo Li, Ziqi Huang, et al. Demystifying video reasoning. *arXiv preprint*
505 *arXiv:2603.16870*, 2026.
- 506 [39] Thaddäus Wiedemer, Yuxuan Li, Paul Vicol, Shixiang Shane Gu, Nick Matarese, Kevin Swersky,
507 Been Kim, Priyank Jaini, and Robert Geirhos. Video models are zero-shot learners and reasoners.
508 *arXiv preprint arXiv:2509.20328*, 2025.
- 509 [40] Bing Wu, Chang Zou, Changlin Li, Duojuan Huang, Fang Yang, Hao Tan, Jack Peng, Jianbing
510 Wu, Jiangfeng Xiong, Jie Jiang, et al. Hunyuanvideo 1.5 technical report. *arXiv preprint*
511 *arXiv:2511.18870*, 2025.
- 512 [41] Siye Wu, Jian Xie, Yikai Zhang, Aili Chen, Kai Zhang, Yu Su, and Yanghua Xiao. Arm:
513 Adaptive reasoning model. *arXiv preprint arXiv:2505.20258*, 2025.
- 514 [42] Zeyue Xue, Siming Fu, Jie Huang, Shuai Lu, Haoran Li, Yijun Liu, Yuming Li, Xiaoxuan
515 He, Mengzhao Chen, Haoyang Huang, et al. A systematic post-train framework for video
516 generation. *arXiv preprint arXiv:2604.25427*, 2026.
- 517 [43] Zeyue Xue, Jie Wu, Yu Gao, Fangyuan Kong, Lingting Zhu, Mengzhao Chen, Zhiheng Liu,
518 Wei Liu, Qiushan Guo, Weilin Huang, et al. Dancegrp: Unleashing grp on visual generation.
519 *arXiv preprint arXiv:2505.07818*, 2025.
- 520 [44] Cheng Yang, Haiyuan Wan, Yiran Peng, Xin Cheng, Zhaoyang Yu, Jiayi Zhang, Junchi Yu,
521 Xinlei Yu, Xiawu Zheng, Dongzhan Zhou, et al. Reasoning via video: The first evaluation of
522 video models’ reasoning abilities through maze-solving tasks. *arXiv preprint arXiv:2511.15065*,
523 2025.
- 524 [45] Ruihan Yang, Prakhar Srivastava, and Stephan Mandt. Diffusion probabilistic modeling for
525 video generation. *Entropy*, 25(10):1469, 2023.
- 526 [46] Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He.
527 Simplerl-zoo: Investigating and taming zero reinforcement learning for open base models in the
528 wild. *arXiv preprint arXiv:2503.18892*, 2025.
- 529 [47] Chenyu Zhang, Daniil Cherniavskii, Antonios Tragoudaras, Antonios Vozikis, Thijmen Nijdam,
530 Derck WE Prinzhorn, Mark Bodracska, Nicu Sebe, Andrii Zadaianchuk, and Efstratios Gavves.
531 Morpheus: Benchmarking physical reasoning of video generative models with real physical
532 experiments. *arXiv preprint arXiv:2504.02918*, 2025.
- 533 [48] Tinghui Zhu, Kai Zhang, Jian Xie, and Yu Su. Deductive beam search: Decoding deducible
534 rationale for chain-of-thought reasoning. *arXiv preprint arXiv:2401.17686*, 2024.

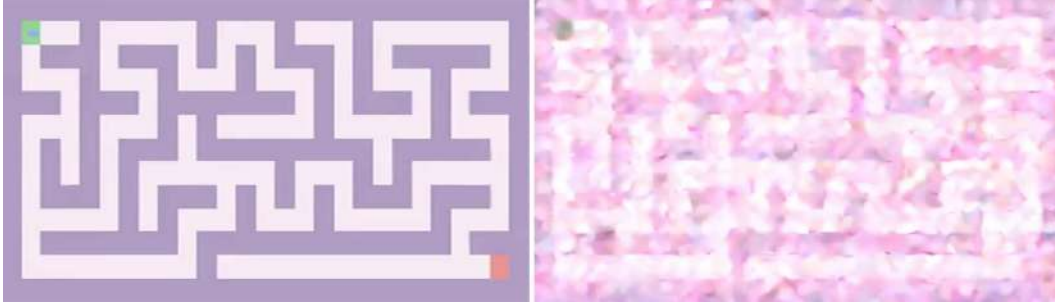


Figure 5: Qualitative example of reward hacking in the absence of KL-divergence regularization. It preserves the wall constraint, and saturates all paths to connect two endpoints, thereby achieving a maximal reward.

535 A Experimental Setup

536 A.1 Dataset

- 537 • **Maze:** We generate 10,000 samples with grid dimensions ranging from 7×7 to 21×21 . Each
538 instance pairs an unsolved layout containing start and goal markers with a ground-truth video that
539 renders the contiguous path between them.
- 540 • **FlowFree:** We generate 10,000 puzzles with grid sizes between 5×5 and 8×8 by splitting
541 Hamiltonian paths into colored segments, ensuring that a valid solution must occupy all available
542 cells. The initial frame displays only the discrete color-pair endpoints, and the video unrolls the
543 progressive coloring of each path.
- 544 • **Sokoban:** This domain includes 10,000 puzzles with grid sizes from 6×6 to 10×10 and 1–3
545 boxes. Solution trajectories are capped at 60 moves. The input frame depicts the initial board
546 configuration, while the video unrolls the solution at a resolution of one agent push per frame,
547 ensuring strict alignment between temporal and logical steps.

548 A.2 Evaluation

549 **Trajectory Alignment Metrics.** To measure the alignment with the ground-truth reference, we
550 compute precision, recall, and F1 at the unit most natural to each task’s solution manifold:

- 551 • **Maze (Pixel-level):** We compute a change mask between the initial and final frames to
552 isolate the generated path from the static background. This ensures the metric captures the
553 model’s intervention rather than background reconstruction.
- 554 • **FlowFree (Cell-level):** We extract mean colors for each grid cell in the terminal frame. This
555 avoids penalizing anti-aliasing artifacts and focuses on the semantic correctness of the path
556 coloring.
- 557 • **Sokoban (Action-level):** Since multiple trajectories can yield the same final state, we decode
558 the video into a symbolic action sequence $a \in \{U(\text{Up}), D(\text{Down}), L(\text{Left}), R(\text{Right})\}$. We
559 report position-aligned F1, which penalizes out-of-order or invalid moves.

560 **Symbolic Success Rate.** Alignment with a single GT reference is insufficient to detect valid
561 alternative solutions or visually plausible but rule-violating outputs. We therefore implement binary
562 success detectors that parse the video \mathbf{V} into symbolic states to verify task-specific rules:

- 563 • **Maze Success:** Requires a connected path between markers without violating wall con-
564 straints.
- 565 • **FlowFree Success:** Validates endpoint preservation, color connectivity, and the grid fill-rate.
- 566 • **Sokoban Success:** Evaluates *process validity* over all frames, checking that player and box
567 displacements follow physics-based rules and the final state matches the target.

568 **B Analysis**

569 **B.1 KL Constraint**

570 **NeurIPS Paper Checklist**

571 The checklist is designed to encourage best practices for responsible machine learning research,
572 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
573 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
574 follow the references and follow the (optional) supplemental material. The checklist does NOT count
575 towards the page limit.

576 Please read the checklist guidelines carefully for information on how to answer these questions. For
577 each question in the checklist:

- 578 • You should answer [Yes], [No], or [N/A].
- 579 • [N/A] means either that the question is Not Applicable for that particular paper or the
580 relevant information is Not Available.
- 581 • Please provide a short (1–2 sentence) justification right after your answer (even for [N/A]).

582 **The checklist answers are an integral part of your paper submission.** They are visible to the
583 reviewers, area chairs, senior area chairs, and ethics reviewers. You will also be asked to include it
584 (after eventual revisions) with the final version of your paper, and its final version will be published
585 with the paper.

586 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
587 While [Yes] is generally preferable to [No], it is perfectly acceptable to answer [No] provided a
588 proper justification is given (e.g., error bars are not reported because it would be too computationally
589 expensive” or “we were unable to find the license for the dataset we used”). In general, answering
590 [No] or [N/A] is not grounds for rejection. While the questions are phrased in a binary way, we
591 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
592 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
593 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
594 please point to the section(s) where related material for the question can be found.

595 **IMPORTANT, please:**

- 596 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**
- 597 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 598 • **Do not modify the questions and only use the provided macros for your answers.**

599 **1. Claims**

600 Question: Do the main claims made in the abstract and introduction accurately reflect the
601 paper’s contributions and scope?

602 Answer: [Yes]

603 Justification: The claims made should match theoretical and experimental results, and reflect
604 how much the results can be expected to generalize to other settings.

605 Guidelines:

- 606 • The answer [N/A] means that the abstract and introduction do not include the claims
607 made in the paper.
- 608 • The abstract and/or introduction should clearly state the claims made, including the
609 contributions made in the paper and important assumptions and limitations. A [No] or
610 [N/A] answer to this question will not be perceived well by the reviewers.
- 611 • The claims made should match theoretical and experimental results, and reflect how
612 much the results can be expected to generalize to other settings.
- 613 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
614 are not attained by the paper.

615 **2. Limitations**

616 Question: Does the paper discuss the limitations of the work performed by the authors?

617 Answer: [TODO]

618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669

Justification: **[TODO]**

Guidelines:

- The answer **[N/A]** means that the paper has no limitation while the answer **[No]** means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate “Limitations” section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren’t acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: **[N/A]**

Justification: The paper does not include theoretical results.

Guidelines:

- The answer **[N/A]** means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: **[Yes]**

Justification: The full experimental setup is presented in Section 5 and Section A.

Guidelines:

- The answer **[N/A]** means that the paper does not include experiments.

- 670 • If the paper includes experiments, a [No] answer to this question will not be perceived
671 well by the reviewers: Making the paper reproducible is important, regardless of
672 whether the code and data are provided or not.
- 673 • If the contribution is a dataset and/or model, the authors should describe the steps taken
674 to make their results reproducible or verifiable.
- 675 • Depending on the contribution, reproducibility can be accomplished in various ways.
676 For example, if the contribution is a novel architecture, describing the architecture fully
677 might suffice, or if the contribution is a specific model and empirical evaluation, it may
678 be necessary to either make it possible for others to replicate the model with the same
679 dataset, or provide access to the model. In general, releasing code and data is often
680 one good way to accomplish this, but reproducibility can also be provided via detailed
681 instructions for how to replicate the results, access to a hosted model (e.g., in the case
682 of a large language model), releasing of a model checkpoint, or other means that are
683 appropriate to the research performed.
- 684 • While NeurIPS does not require releasing code, the conference does require all submis-
685 sions to provide some reasonable avenue for reproducibility, which may depend on the
686 nature of the contribution. For example
 - 687 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
688 to reproduce that algorithm.
 - 689 (b) If the contribution is primarily a new model architecture, the paper should describe
690 the architecture clearly and fully.
 - 691 (c) If the contribution is a new model (e.g., a large language model), then there should
692 either be a way to access this model for reproducing the results or a way to reproduce
693 the model (e.g., with an open-source dataset or instructions for how to construct
694 the dataset).
 - 695 (d) We recognize that reproducibility may be tricky in some cases, in which case
696 authors are welcome to describe the particular way they provide for reproducibility.
697 In the case of closed-source models, it may be that access to the model is limited in
698 some way (e.g., to registered users), but it should be possible for other researchers
699 to have some path to reproducing or verifying the results.

700 5. Open access to data and code

701 Question: Does the paper provide open access to the data and code, with sufficient instruc-
702 tions to faithfully reproduce the main experimental results, as described in supplemental
703 material?

704 Answer: [TODO]

705 Justification: [TODO]

706 Guidelines:

- 707 • The answer [N/A] means that paper does not include experiments requiring code.
- 708 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
709 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 710 • While we encourage the release of code and data, we understand that this might not
711 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
712 including code, unless this is central to the contribution (e.g., for a new open-source
713 benchmark).
- 714 • The instructions should contain the exact command and environment needed to run to
715 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
716 neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 717 • The authors should provide instructions on data access and preparation, including how
718 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 719 • The authors should provide scripts to reproduce all experimental results for the new
720 proposed method and baselines. If only a subset of experiments are reproducible, they
721 should state which ones are omitted from the script and why.
- 722 • At submission time, to preserve anonymity, the authors should release anonymized
723 versions (if applicable).

- 724 • Providing as much information as possible in supplemental material (appended to the
725 paper) is recommended, but including URLs to data and code is permitted.

726 **6. Experimental setting/details**

727 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
728 rameters, how they were chosen, type of optimizer) necessary to understand the results?

729 Answer: [Yes]

730 Justification: The full experimental setup is presented in Section 5 and Section A.

731 Guidelines:

- 732 • The answer [N/A] means that the paper does not include experiments.
733 • The experimental setting should be presented in the core of the paper to a level of detail
734 that is necessary to appreciate the results and make sense of them.
735 • The full details can be provided either with the code, in appendix, or as supplemental
736 material.

737 **7. Experiment statistical significance**

738 Question: Does the paper report error bars suitably and correctly defined or other appropriate
739 information about the statistical significance of the experiments?

740 Answer: [No]

741 Justification: Error bars are not reported because it would be too computationally expensive.
742 Inference one time would cost more than 24 GPU hours of B200.

743 Guidelines:

- 744 • The answer [N/A] means that the paper does not include experiments.
745 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
746 intervals, or statistical significance tests, at least for the experiments that support the
747 main claims of the paper.
748 • The factors of variability that the error bars are capturing should be clearly stated (for
749 example, train/test split, initialization, random drawing of some parameter, or overall
750 run with given experimental conditions).
751 • The method for calculating the error bars should be explained (closed form formula,
752 call to a library function, bootstrap, etc.)
753 • The assumptions made should be given (e.g., Normally distributed errors).
754 • It should be clear whether the error bar is the standard deviation or the standard error
755 of the mean.
756 • It is OK to report 1-sigma error bars, but one should state it. The authors should
757 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
758 of Normality of errors is not verified.
759 • For asymmetric distributions, the authors should be careful not to show in tables or
760 figures symmetric error bars that would yield results that are out of range (e.g., negative
761 error rates).
762 • If error bars are reported in tables or plots, the authors should explain in the text how
763 they were calculated and reference the corresponding figures or tables in the text.

764 **8. Experiments compute resources**

765 Question: For each experiment, does the paper provide sufficient information on the com-
766 puter resources (type of compute workers, memory, time of execution) needed to reproduce
767 the experiments?

768 Answer: [Yes]

769 Justification: The full experimental setup is presented in Section 5 and Section A.

770 Guidelines:

- 771 • The answer [N/A] means that the paper does not include experiments.
772 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
773 or cloud provider, including relevant memory and storage.

- 774
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- 775
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 776
- 777
- 778

779 9. Code of ethics

780 Question: Does the research conducted in the paper conform, in every respect, with the
781 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

782 Answer: [Yes]

783 Justification: The research conducted in the paper conform, in every respect, with the
784 NeurIPS Code of Ethics.

785 Guidelines:

- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 786
- 787
- 788
- 789
- 790
- 791

792 10. Broader impacts

793 Question: Does the paper discuss both potential positive societal impacts and negative
794 societal impacts of the work performed?

795 Answer: [TODO]

796 Justification: [TODO]

797 Guidelines:

- The answer [N/A] means that there is no societal impact of the work performed.
 - If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).
- 798
- 799
- 800
- 801
- 802
- 803
- 804
- 805
- 806
- 807
- 808
- 809
- 810
- 811
- 812
- 813
- 814
- 815
- 816
- 817
- 818
- 819

820 11. Safeguards

821 Question: Does the paper describe safeguards that have been put in place for responsible
822 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
823 image generators, or scraped datasets)?

824 Answer: [N/A]

825 Justification: This work does not have such risks.

826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876

Guidelines:

- The answer [N/A] means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [N/A]

Justification: The paper does not use existing assets.

Guidelines:

- The answer [N/A] means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [TODO]

Justification: [TODO]

Guidelines:

- The answer [N/A] means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919

Answer: [N/A]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [N/A]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does *not* impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [N/A]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

Guidelines:

- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.