

長尺動画生成タスクにおける メタ評価ベンチマーク



C2-01

松田 陵佑¹, 工藤 慧音^{1,2}, 吉田 遥音¹, 清水 伸幸³, 鈴木 潤^{1,2,4}
 (¹ 東北大学, ² 理化学研究所, ³ LINEヤフー株式会社, ⁴ 国立情報科学研究所 LLC)
 is-failab-research@grp.tohoku.ac.jp

概要

- 長尺動画生成モデルの評価システムの性能をメタ評価するためのSLVMEval ベンチマークを提案
- 人工的に平均19分 最大約3時間の高品質/低品質の長尺動画ペアを作成
- 既存の評価システムは9観点で人間性能に及ばない
特にプロンプトと動画の一貫性に関して性能が低い

背景

- 動画生成(T2V)モデルの性能は進化しているが長尺動画生成(T2LV)と**その評価手法は未開拓**
- ❓ 正しく測定できないものは改善出来ないのでは?
- 既存のT2V評価モデルは数秒の動画で学習された評価モデル [He+’2024]
 - 人間が評価するのはコスト・時間的に困難
- ❗ 評価モデルをメタ評価するベンチマークの整備が必要

データセット

SLVMEval の概要

現在のT2Vモデルの生成時間を大きく超える&観点別の精緻な評価を可能にしたい!
 → 長尺動画に評価観点別に人工的に劣化処理を施す



劣化処理と動画例

見た目の品質



プロンプトと動画の一貫性



実験設定

設定1.ペアワイズ評価

元動画 vs. 劣化動画 のペアのどちらが高品質かを評価システムに評価させて正答率を測る

Q. どちらの動画がある観点において高品質か?



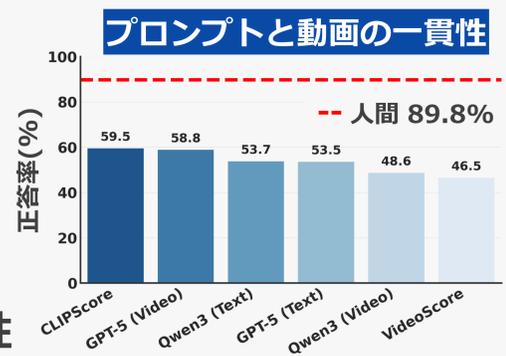
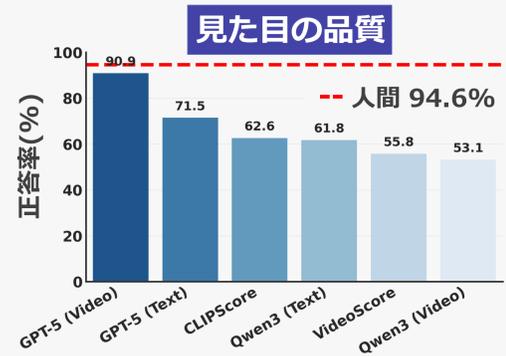
設定2.評価システム

- VLM-as-a-judge** (GPT-5, Qwen3-VL)
 - 動画ベース評価: VLMに2つの動画を入力し評価
 - テキストベース評価: VLMに事前にキャプションを作成させて2つのキャプションを入力し評価
- CLIPScore** [Hessel+’2021]
- VideoScore** [He+’2024]

実験結果

結果1.全体の分析

- 人間は**90%以上**の正答率
- ❗ **GPT-5でも9観点で人間に及ばない**
- 既存のVLM等は簡単なタスクでさえ解けない
- 特にプロンプトと動画の一貫性は性能が低い



結果2.動画長への脆弱性

- 人間は**頑健性が高い**
- ❗ 評価システムは複数観点で動画長の増加に伴い正答率が低下の傾向

