# VIEScore: Towards Explainable Metrics for Conditional Image Synthesis Evaluation

♠Max Ku, ♠Dongfu Jiang, ♠Cong Wei, ♥Xiang Yue, ♠Wenhu Chen

♠University of Waterloo, ♥IN.AI Research
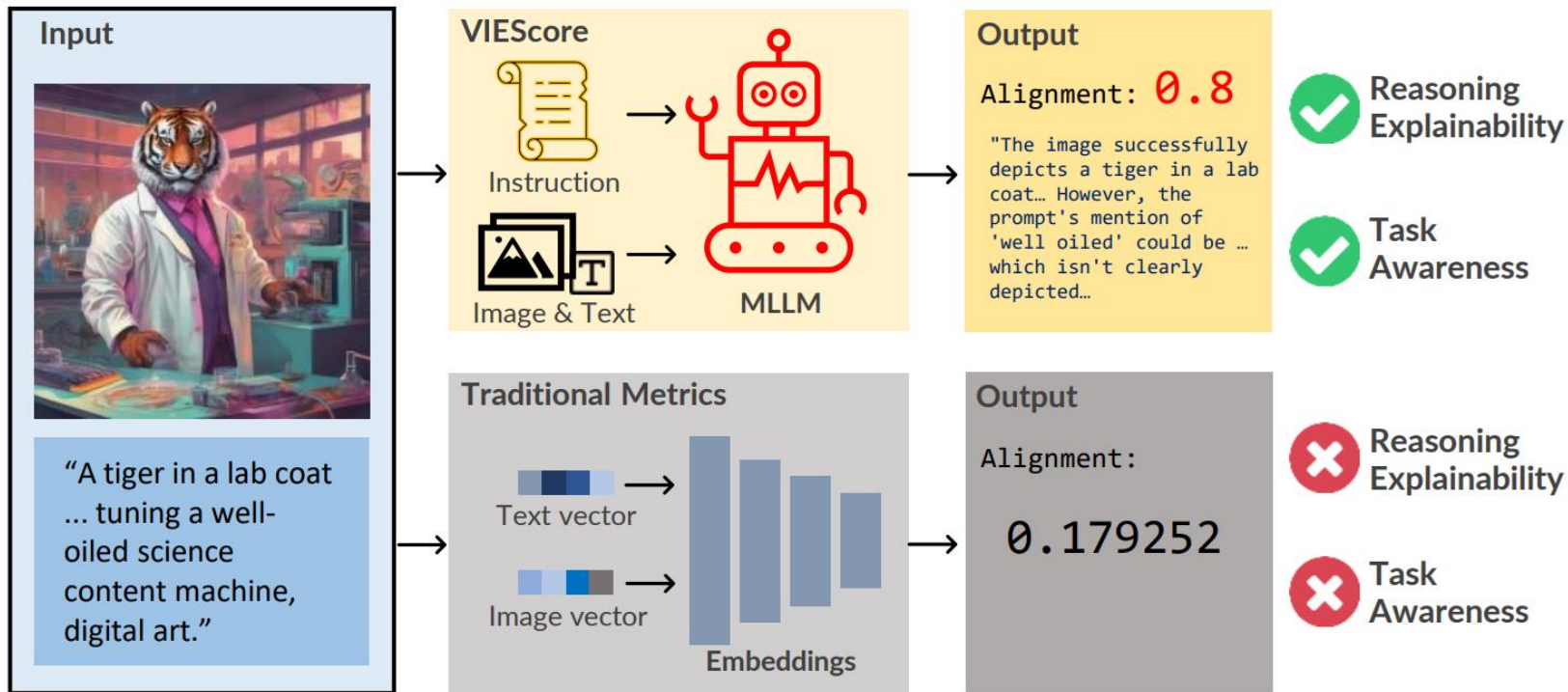
tiger-ai-lab.github.io/VIEScore

Paper ID: #2258

ACL 2024
Bangkok, Thailand

UNIVERSITY OF WATERLOO

# Motivation

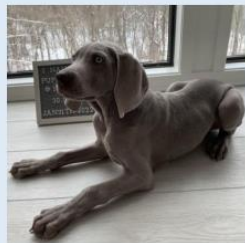*Metrics in the future would provide not just the score but also the rationale*

- Understanding of judgment per instance is needed

- Traditional metrics are not task aware

# Visual Instruction-guided Explainable Score (VIEScore)

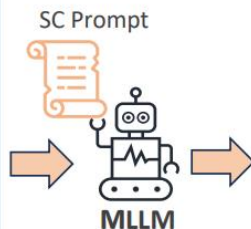# How VIEScore works?



**Input**

Concept 1   Concept 2

"dog sitting in a driving car"

Prompt for Synthesized Image

Synthesized Image

SC Prompt

MLLM

PQ Prompt

MLLM

**Semantic Consistency (SC)**

Response:
The dog isn't sitting as a driver would, hence the score of **7** for **following the prompt**. The dog in the second image strongly resembles..., warranting a score of **9** for **resemblance**. The car's interior and style are entirely different, which results in a score of **0** for **resemblance**...
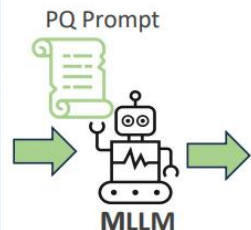
SC Scores:
Alignment with the prompt: **7**
Resemblance to concept 1: **9**
Resemblance to concept 2: **0**

SC score = min(7, 9, 0) = **0**

**Perceptual Quality (PQ)**

Response:
The **naturalness score** is given a **7** because the dog appears well integrated into the car setting with proper shadowing and lighting that matches the interior of the car...The **artifact score** is an **8** because the image is clear ...

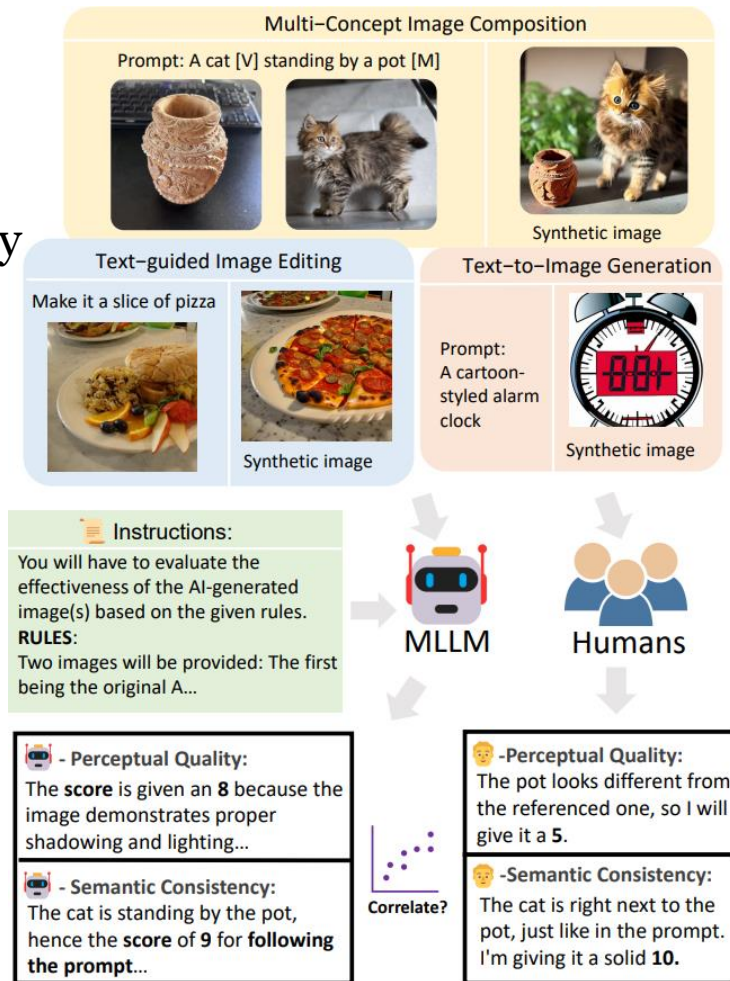PQ Scores:
looks natural: **7**
Has no artifacts: **8**

PQ score = min(7, 8) = **7**

# Experiment Setup

A wide range of image synthesis tasks study

- Correlation of VIEScore to Human

- V.S.

- Correlation of Traditional metrics to Human

*Where can we get this kind of human annotation data?*

# Experiment Setup (Cont.)

- Human data from ImagenHub(ICLR 2024)

- 29 Models across 7 tasks
  - Total 14403 annoations

- Each annotation 3 human metrics:
  - SC : Conditions-Image alignment
  - PQ: Realism and Natural sense
  - Overall: sqrt(SC x PQ)

- Human guideline is used as prompt

ImagenHub: Standardizing the evaluation of conditional image generation models
https://arxiv.org/abs/2310.01596

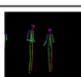| $c_1$ | $c_2$ | $c_2$ | Task | $y$ |
|---|---|---|---|---|
| A cartoon styled alarm clock | Ø | Ø | Text-to-Image Generation | |
| | | Change frisbee to a football | Mask-guided Image Editing | |
| | Make it a slice of pizza instead of the sandwich | Ø | Text-guided Image Editing | |
| | A [V] dog in the Versailles hall of mirrors | Ø | Subject-Driven Image Generation | |
| | | Replace glasses with [V] glasses | Subject-Driven Image Editing | |
| | | A cat [V] standing by a pot [M] | Multi-Concept Image Composition | |
| | A small dog is curled up on top of the shoes | Ø | Control-guided Image Generation | |

# Main Result



Performance Across 7 Tasks

# Why one-shot setting achieve worse performance?

- MLLMs struggle in In-Context Learning when multiple images exists
  - Reasoning is affected

- Appears on all MLLMs we benchmarked



**Prompt**

....... (Detailed text of rating instruction on PQ) .......

**1st image as a rating example.**
PQ scores:
Image looks natural? 5
Image has no artifacts? 5
Reasoning:
The image gives an unnatural feeling on hands of the girl. There is also minor distortion on the eyes of the girl.
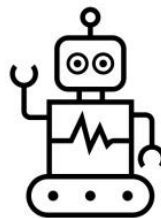
**Please evaluate the 2nd image.**
PQ scores:
Image looks natural? _
Image has no artifacts? _
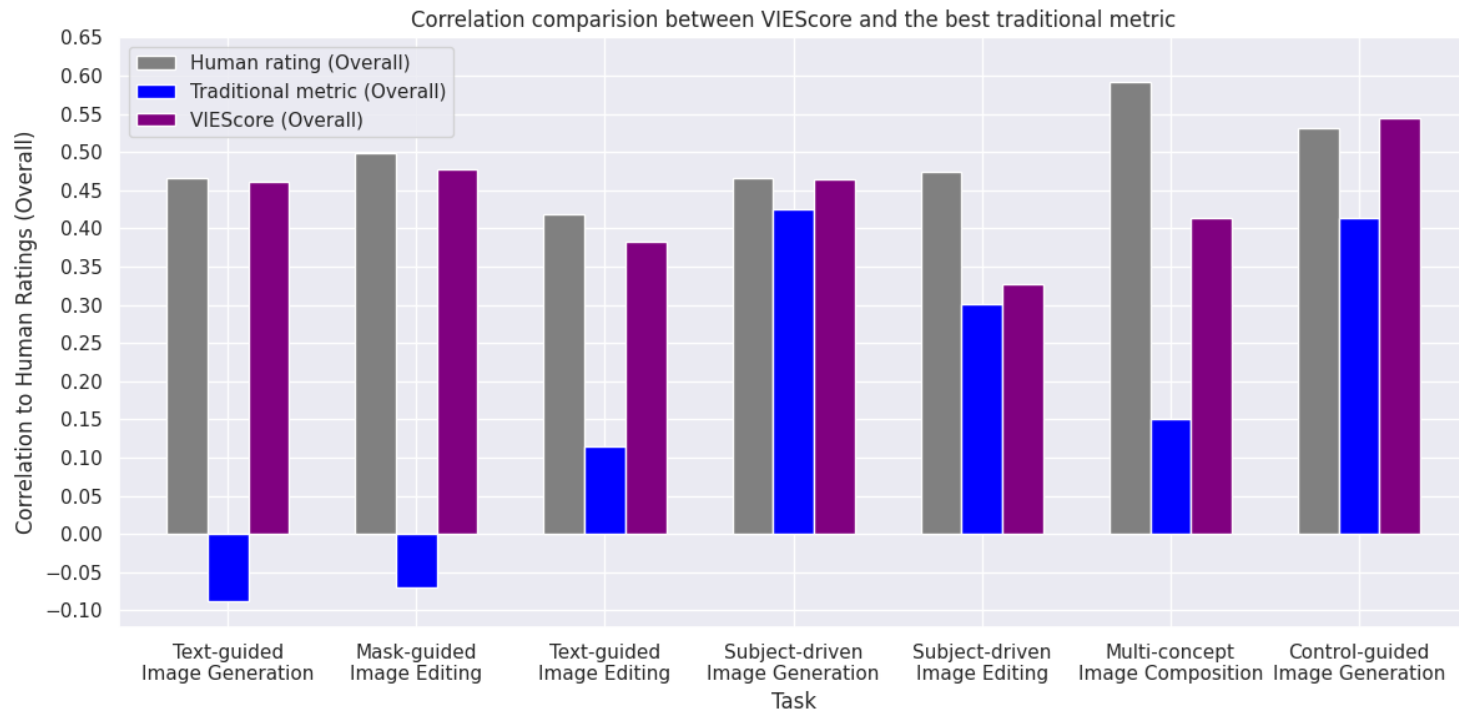Reasoning:
_____
_____
_____
_____

**Response**

PQ scores:
Image looks natural? 3
Image has no artifacts? 4
Reasoning:
The girl's image has an unnatural blurring effect .... The birds also look slightly distorted. The cat's the face looks slightly artificial.

LLM

# VIEScore achieve higher correlation than Traditional metrics



Correlation comparision between VIEScore and the best traditional metric

# There's still obstacles toward explainable metrics

- MLLMs struggles to see minor difference when given 2 similar images

# Code Release

- Code available on GitHub

- Easy to use and modify

```python
from viescore import VIEScore
backbone = "gemini"
vie_score = VIEScore(backbone=backbone, task="t2v")

score_list = vie_score.evaluate(pil_image, text_prompt)
sementics_score, quality_score, overall_score = score_list
```