Re-Thinking the Automatic Evaluation of Image-Text Alignment in Text-to-Image Models

Anonymous ACL submission

Abstract

Text-to-image models often struggle to generate images that precisely match textual prompts. Prior research has extensively studied the evaluation of image-text alignment in text-to-image generation. However, existing evaluations primarily focus on agreement with human assessments, neglecting other critical properties of a trustworthy evaluation framework. In this work, we first identify two key aspects that a reliable evaluation should address. We then empirically demonstrate that current mainstream evaluation frameworks fail to fully satisfy these properties across a diverse range of metrics and models. Finally, we propose recommendations for improving image-text alignment evaluation.

1 Introduction

027

030

034

036

Text-to-Image (T2I) models have demonstrated remarkable capabilities in generating high-quality, realistic images (Betker et al., 2023; Esser et al., 2024). Despite these advances, they still face challenges in accurately interpreting and adhering to user prompts. Common failures include generating incorrect object counts, attributes, or spatial relationships (Li et al., 2024). Nevertheless, evaluating text-image alignment remains a persistent and unresolved problem in the field.

Several frameworks exist for evaluating imagetext alignment. Model-based approaches include CLIPScore (Hessel et al., 2021) and VQAScore (Lin et al., 2024). Component-based methods decompose text prompts into fine-grained elements and assess alignment through techniques like question generation and answering (QG/A) (Hu et al., 2023; Cho et al., 2023). Additionally, detectorbased frameworks such as UniDet-Eval (Huang et al., 2023) leverage object detection for evaluation.

Despite the vast number of evaluation frameworks, few studies have thoroughly investigated "A sculpture of a green hot rod on a city sidewalk"



VQAScore = 43.85

VQAScore = 33.49

Figure 1: A robustness failure case of VQAScore. The right image is visually the same as the left one, yet their calculated VQAScore differs a lot.

the trustworthiness of image-text alignment assessment. Most existing work focuses narrowly on aligning automatic evaluation results with human judgments (Li et al., 2024; Wiles et al., 2025). When addressing trustworthiness, researchers primarily concentrate on improving human evaluation protocols to better validate automatic metrics (Otani et al., 2023; Wiles et al., 2025). However, these approaches overlook the critical point that a truly trustworthy evaluation framework must encompass more dimensions than mere correlation with human assessment. Additional discussion of related works is provided in Appendix A.

In this work, we identify two critical properties for trustworthy evaluation frameworks: **Robustness** and **Significance**. Through empirical analysis, we demonstrate that current evaluation methods fail to fully satisfy these criteria, highlighting an important research direction for improving text-image alignment assessment. For example, as illustrated in Figure 1, while two visually similar images receive substantially different VQAScore evaluations, revealing clear deficiencies in Robustness.

In summary, our work makes two key contributions: (1) identifying Robustness and Significance as fundamental requirements for trustworthy evaluation frameworks, and (2) systematically demon061

062

063

064

065

066

067

041

069

087

095

096

098

100

101

102

103

104

105

107

109

110

strating how current frameworks fail to meet these criteria.

2 Methodology

2.1 Preliminaries

We would like to briefly introduce the common image-text alignment evaluation framework that we discuss in this research. Generally, given a benchmark consisting a set of N prompts P = $\{p_1, ..., p_N\}$, a text-to-image model M is used to generate images $I = \{I_1, ..., I_N\}$ based on these prompts. ¹ An automatic metric J is used to evaluate the alignment between the image-text pair (p_i, I_i) , and outputs a score $s_i = J(p_i, I_i)$. The final evaluation result is the average of all scores: $s_M = \frac{1}{N} \sum_{i=1}^N s_i$. If there are multiple models $M_1, ..., M_K$, the scores $s_{M_1}, ..., s_{M_K}$ also provides a ranking of these models, which is also a part of the evaluation result that people focus on.

2.2 Aspects

We would like to propose two important aspects that a trustworthy image-text alignment evaluation should focus on.

Robustness Robustness requires that evaluation results remain consistent under reasonable perturbations of the input pair (p_i, I_i) . In this work, we specifically examine two critical dimensions of robustness: (1) robustness to randomness and (2) robustness to image perturbations.

Robustness to Randomness Most state-ofthe-art text-to-image models employ denoising diffusion processes, where the generation output for a given prompt p_i depends on the sampled noise prior. This inherent randomness introduces variability in evaluation results. A trustworthy evaluation framework must maintain consistent model rankings despite this randomness - otherwise, the evaluation fails to reliably indicate which model performs better. To assess robustness under randomness, we systematically evaluate model performance across different random seeds.

Robustness to Image Perturbation We make a fundamental assumption that for visually similar images I_i and I'_i , their evaluation scores $J(p_i, I_i)$ and $J(p_i, I'_i)$ should also be close. Large discrepancies would indicate potential metric flaw rather than true model capability. To assess robustness against image perturbations, we apply a minimal transformation: given an image I with pixel values $I_{c,h,w} \in [0, 255], c, h, w$ corresponds to the channels, height and width of the image, we achieve a perturbed image I' with pixel values ²:

$$I_{c,h,w}^{'} = \begin{cases} I_{c,h,w} + 1 & \text{if } I_{c,h,w} < 255\\ I_{c,h,w} & \text{otherwise} \end{cases}$$
(1)

The robustness metric is computed by calculating the performance gap as:

$$\Delta J_{i} = |J(p_{i}, I_{i}) - J(p_{i}, I_{i}^{'})|$$
(2)

Significance The property of Significance examines whether an observed performance difference (e.g., $s_{M_1} > s_{M_2}$) reflects a meaningful superiority of model M_1 over M_2 . To quantify this, we employ two complementary approaches:

- 1. Statistical Testing: For evaluation score sets $S_{M_1} = (s_{M_1}^1, ..., s_{M_1}^N)$ and $S_{M_2} = (s_{M_2}^1, ..., s_{M_2}^N)$, we conduct a paired t-test to assess statistical significance.
- 2. Dominance Ratio: We compute the empirical probability of M_1 over M_2 as:

$$R = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}[s_{M_1}^i > s_{M_2}^i]$$
(3)

3 Experiment Setup

For our evaluation, we select four widely-used textto-image generation models that represent diverse architectures and capabilities: Stable-Diffusion-3 (SD3) (Esser et al., 2024), Stable-Diffusion-XL (SD-XL) (Podell et al., 2023), Stable-Diffusion-1.5 (SD1.5) (Rombach et al., 2022), and PixArt-Sigma-XL (Pixart) (Chen et al., 2024a). This carefully chosen set ensures comprehensive coverage of the current state-of-the-art in diffusion-based text-toimage generation.

For metrics, we employ three widely applicable metrics: CLIPScore (Hessel et al., 2021), VQAScore (Li et al., 2024), and DSGScore (Cho et al., 2023). For benchmarking, we use MSCOCO (Lin 128

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

129 130

131 132

133

134

135

136

138

139

140

141

142

143

144

145

146

147

148

¹Though there can be multiple images generated using a certain prompt, many practices only generate one image per prompt.

²We confirm the perturbation's visual imperceptibility through manual inspection of multiple cases, verifying that modified images remain indistinguishable from originals.

Model Name	VQAScore			CLIPScore			DSGScore		
	42	3407	5096	42	3407	5096	42	3407	5096
Stable-Diffusion-3	91.18(1)	90.90(1)	91.06(1)	26.39(1)	26.34(1)	26.36(1)	93.66(1)	91.99(1)	93.52(1)
Stable-Diffusion-XL	86.63(3)	86.01(3)	85.77(3)	25.93(2)	25.81(2)	25.85(2)	89.68(3)	90.04(3)	90.44(2)
Pixart-Sigma-XL	87.04(2)	87.16(2)	86.72(2)	25.78(3)	25.71(3)	25.75(4)	90.52(2)	90.53(2)	89.99(3)
Stable-Diffusion-1.5	76.26(4)	75.79(4)	77.32(4)	25.76(4)	25.58(4)	25.76(3)	83.23(4)	82.64(4)	83.88(4)

Table 1: Results on MSCOCO using different metrics and different random seeds. The value under metric name is the corresponding random seed used. The value in the brackets is the ranking under the same seed.

et al., 2014), selecting 1,000 prompts from its validation split to assess general text-to-image generation capabilities following common practice (Esser et al., 2024). Additional implementation details are provided in Appendix B.

4 Experiment Results and Analysis

4.1 Analysis of Robustness

150

151

152

153

154

155

156

157

158

159

160

162

163

164

165

166

168

169

170

171

172

173

174

175

177

178

179

180

181

182

186

189

Robustness to Randomness Our investigation first examines evaluation robustness under random seed variations. Table 1 reveals that both CLIP-Score and DSGScore produce inconsistent model rankings across different random seeds, demonstrating their failure to maintain robust evaluation outcomes.

Notably, this inconsistency cannot be solely attributed to narrow score margins between models. For instance, CLIPScore with seed 3407 shows a 0.10 performance gap between Pixart and SD-XL, compared to a larger 0.13 gap between Pixart and SD-1.5. Despite this greater margin, CLIPScore inconsistently ranks Pixart versus SD-1.5 while maintaining stable rankings for SD-XL.

We emphasize that our analysis does not presuppose the rankings produced by VQAScore are inherently "correct." Rather, our core argument establishes that any metric failing to maintain consistent rankings under random seed variations cannot be considered trustworthy, as there exists no reliable basis to determine which ranking might be "correct" when results fluctuate.

Furthermore, while 1000 prompts constitute a large-scale evaluation in text-to-image research, our findings reveal significant robustness failures even at this scale. This persistent inconsistency underscores fundamental challenges in current evaluation practices.

Takeaway 1: An evaluation failing to provide robust ranking under randomness should be viewed less trustworthy, like CLIPScore and DSGScore. **Robustness to Image Perturbation** For our image perturbation analysis, we utilize SD-3 (selected for its superior generation performance). The evaluation employs a fixed random seed of 42, with complete results presented in Table 2.

190

191

192

193

194

195

196

197

198

199

200

201

203

204

205

207

208

209

210

211

212

213

214

215

216

217

218

219

221

222

223

	Avg. ΔJ_i	Max. ΔJ_i
VQAScore CLIPScore	0.44	10.36 7.30
CLIPScore DSGScore	0.74 3.09	7.30 50.00

Table 2: The average and maximum performance gap between original image and perturbated image of different metrics.

The results reveal that even a minimal perturbation of 1 pixel value creates significant performance variations. While VQAScore demonstrates relatively better robustness with an average performance gap of 0.44, both CLIPScore and DSGScore exhibit unacceptably large variations, indicating fundamental robustness limitations.

More concerning are the worst-case scenarios, where this simple perturbation produces dramatic performance gaps across all metrics. Figure 1 illustrates one such failure case for VQAScore. These extreme cases are particularly problematic as they not only occur unpredictably, but also reveal potential for metric exploitation. Further, this worse case may confuse model development when visually identical inputs produce substantially different evaluations.

Takeaway 2: All three metrics fail to maintain a worst case robustness under a simple slight perturbation of image. CLIPScore and DSGScore even fail on average case, revealing a worrying fact of the trustworthiness of their evaluation result.

4.2 Analysis of Significance

We present the p-value of paired T-test and dominance ratio R to explore the significance of the comparison using different metrics in Figure 2.

We first examine the statistical significance of model comparisons. Following conventional stan-



Figure 2: The T-test p-value and better ratio R between models. The value at *i*-th row and *j*-th column in a matrix represents the result of model *i* compared against model *j*. The random seed used is shown in the title of the corresponding heatmap.

dards, we consider results with p-value < 0.05 as statistically significant. Our analysis using VQAScore reveals several notable findings: SD-3 demonstrates statistically significant superiority over SD-XL, Pixart, and SD-1.5. More surprisingly, Pixart shows significant improvement over SD-XL (87.16 vs 86.01) in Table 1, seed 3407), indicating that even a score difference of 1 can reflect meaningful performance gaps between models.

224

225

227

234

240

241

242

243

246

249

252

This phenomenon is not unique to VQAScore. CLIPScore exhibits a similar pattern, showing SD-XL significantly superior to Pixart (25.93 vs 25.78). Similarly, DSGScore suggests potential superiority of Pixart over SD-XL (p = 0.10), despite their small score difference (90.52 vs 89.68). These consistent findings across metrics challenge conventional assumptions, suggesting that statistically significant improvements may occur with smaller metric differences than previously believed. So are current standards for determining meaningful model improvements excessively stringent?

To further investigate this phenomenon, we examine the dominance ratio - the probability that model i generates superior results to model j for a given prompt. We reach another shocking observation that, even if model i is significantly better than model j, the generation result of model i may not bear a large probability of being better than model j. Considering the significance derived using VQAS-

core as mentioned before, actually SD-XL bears a 50% probability generating a "better" result, while Pixart bears 49%, even less than SD-XL! Further, SD-3 just bears 60% probability of generating better results, while SD-XL still bears 40% probability of generating better results, even with a VQAScore gap of 4.5.

253

254

255

256

257

258

259

260

261

262

263

264

265

268

269

270

271

272

273

274

275

276

277

278

280

In this context, the problem is not about metrics only. The problem is, how we view "significance". If we simply view "significance" as a statistical test, we can happily accept a small improvement in metrics since it is enough to demonstrate this significance. However, if we would like to guarantee a better generation performance, it is still essential to seek an even larger metric improvement.

Takeaway 3: A small difference between metrics is enough to reveal a "significance" in statistical analysis. However, this "significance" may not be directly interpreted to actual model performance, while a trustworthy evaluation should take both aspects into account.

5 Conclusion

In this work, we introduce two important aspects, Robustness and Significance, that evaluation frameworks should focus. We conduct a wide range of experiments and reveal some important conclusions that should be taken into consideration when conducting future evaluation on image-text alignment. 281

Limitations

References

2(3):8.

The main limitation of this work is that it does not

provide a better evaluation framework to address

the problems discussed in this paper. A better eval-

James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jian-

feng Wang, Linjie Li, Long Ouyang, Juntang Zhuang,

Joyce Lee, Yufei Guo, and 1 others. 2023. Improving

image generation with better captions. Computer Sci-

ence. https://cdn. openai. com/papers/dall-e-3. pdf,

Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu,

Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping

Luo, Huchuan Lu, and Zhenguo Li. 2024a. Pixart- σ :

Weak-to-strong training of diffusion transformer for

Xiang Chen, Chenxi Wang, Yida Xue, Ningyu Zhang, Xiaoyan Yang, Qiang Li, Yue Shen, Lei Liang, Jinjie

Gu, and Huajun Chen. 2024b. Unified hallucina-

tion detection for multimodal large language models.

Jaemin Cho, Yushi Hu, Roopal Garg, Peter Ander-

son, Ranjay Krishna, Jason Baldridge, Mohit Bansal,

Jordi Pont-Tuset, and Su Wang. 2023. Davidsonian

scene graph: Improving reliability in fine-grained

evaluation for text-image generation. arXiv preprint

Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim

Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, and 1

others. 2024. Scaling rectified flow transformers for

high-resolution image synthesis. In Forty-first Inter-

Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig

Schmidt. 2023. Geneval: An object-focused frame-

work for evaluating text-to-image alignment. In Ad-

vances in Neural Information Processing Systems,

volume 36, pages 52132–52152. Curran Associates,

Tejas Gokhale, Hamid Palangi, Besmira Nushi, Vib-

hav Vineet, Eric Horvitz, Ece Kamar, Chitta Baral,

and Yezhou Yang. 2023. Benchmarking spatial re-

lationships in text-to-image generation. Preprint,

Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. Clipscore: A

reference-free evaluation metric for image captioning.

In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages

Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang,

Mari Ostendorf, Ranjay Krishna, and Noah A. Smith.

national Conference on Machine Learning.

4k text-to-image generation. arXiv.

arXiv preprint arXiv:2402.03190.

arXiv:2310.18235.

arXiv:2212.10015.

7514-7528.

Inc.

uation framework is left for future works.

20

28

28

286 287 288

290 291 292

293

299 300 301

303 304 305

3

310

316 317 318

319 320

321 322

322

324 325

326

3

3

330 331

332

333

2023. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering. In *Proceedings of the IEEE/CVF International Con-ference on Computer Vision (ICCV)*, pages 20406–20417.

334

335

337

338

339

340

341

342

343

344

345

348

350

351

353

354

355

356

358

359

360

361

362

363

364

365

366

367

370

371

372

373

374

375

376

377

379

381

384

385

388

389

- Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. 2023. T2i-compbench: A comprehensive benchmark for open-world compositional textto-image generation. *Preprint*, arXiv:2307.06350.
- Baiqi Li, Zhiqiu Lin, Deepak Pathak, Jiayao Li, Yixin Fei, Kewen Wu, Xide Xia, Pengchuan Zhang, Graham Neubig, and Deva Ramanan. 2024. Evaluating and improving compositional text-to-visual generation. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition (CVPR) Workshops, pages 5290–5301.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pretraining for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In Computer Vision– ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13, pages 740–755. Springer.
- Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. 2024. Evaluating text-to-visual generation with image-to-text generation. *arXiv preprint arXiv:2404.01291*.
- Mayu Otani, Riku Togashi, Yu Sawai, Ryosuke Ishigami, Yuta Nakashima, Esa Rahtu, Janne Heikkilä, and Shin'ichi Satoh. 2023. Toward verifiable and reproducible human evaluation for text-toimage generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14277–14286.
- Maitreya Patel, Tejas Gokhale, Chitta Baral, and Yezhou Yang. 2024. Conceptbed: Evaluating concept learning abilities of text-to-image diffusion models. *Preprint*, arXiv:2306.04695.
- Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. 2023. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.

- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. Highresolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695.
 - Olivia Wiles, Chuhan Zhang, Isabela Albuquerque, Ivana Kajić, Su Wang, Emanuele Bugliarello, Yasumasa Onoe, Pinelopi Papalampidi, Ira Ktena, Chris Knutsen, Cyrus Rashtchian, Anant Nawalgaria, Jordi Pont-Tuset, and Aida Nematzadeh. 2025. Revisiting text-to-image evaluation with gecko: On metrics, prompts, and human ratings. *Preprint*, arXiv:2404.16820.
 - Xindi Wu, Dingli Yu, Yangsibo Huang, Olga Russakovsky, and Sanjeev Arora. 2024. Conceptmix: A compositional image generation benchmark with controllable difficulty. *Preprint*, arXiv:2408.14339.
 - Michal Yarom, Yonatan Bitton, Soravit Changpinyo, Roee Aharoni, Jonathan Herzig, Oran Lang, Eran Ofek, and Idan Szpektor. 2024. What you see is what you read? improving text-image alignment evaluation. Advances in Neural Information Processing Systems, 36.
 - Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.

A Related Works

400

401

402

403

404

405

406

407

408 409

410

411

412 413

414

415

416

417

418

419 420

421

422

423

424 425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

There are many metrics and benchmarks focusing on evaluating image-text alignment in text-toimage generation.

CLIP-Score (Radford et al., 2021; Hessel et al., 2021) evaluates image-text alignment by computing the cosine similarity of CLIP embeddings. VQAScore (Lin et al., 2024) queries a VQA model to determine if the image corresponds to the prompt, using the "Yes" logit as the metric. T2I-CompBench (Huang et al., 2023) leverages MiniGPT-4 (Zhu et al., 2023) CoT to generate an alignment score.

Decomposition-based metrics break down text prompts into smaller components and assess the accuracy of each part. TIFA (Hu et al., 2023) generates visual questions and uses a VQA model to verify the correctness of each component. T2I-CompBench (Huang et al., 2023) employs Blip-VQA (Li et al., 2022), while DavidSceneGraph (Cho et al., 2023) and VQ² (Yarom et al., 2024) are similar to TIFA. MHalu-Bench (Chen et al., 2024b) builds a pipeline of tools to check the correctness of each component directly. ConceptMix(Wu et al., 2024) is a more complicated benchmark. Some other benchmarks includes (Gokhale et al., 2023; Patel et al., 2024), GenEval (Ghosh et al., 2023).

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

However, all these works consider only alignment with human annotation to evaluate the validity of their evaluation. (Wiles et al., 2025) takes a step forward and points out different human evaluation template influences evaluation results, but it still focuses only on human annotation. To the best of our knowledge, we are the first to explore imagetext alignment evaluation from the inner properties of the evaluation.

B Experiment Details

We use the default inference hyper-parameter of each model used. We list the details as follows:

Model Name	T	w
Stable-Diffusion-3	28	7.0
Stable-Diffusion-XL	50	5.0
Stable-Diffusion-1.5	50	7.5
PixArt- Σ -XL	20	4.5

Table 3: Details of our inference hyper-parameter. T represents total denoising steps and w represents guidance scale.