

EvalGIM: A Library for Evaluating Generative Image Models

Melissa Hall^{*1}, Oscar Mañas^{*1,2}, Reyhane Askari-Hemmat^{*1}, Mark Ibrahim^{*1}, Candace Ross^{*1}, Pietro Astolfi^{*1}, Tariq Berrada Ifriqi^{*1,3}, Marton Havasi^{*1}, Yohann Benchetrit¹, Karen Ullrich¹, Carolina Braga¹, Abhishek Charnalia¹, Maeve Ryan¹, Mike Rabbat¹, Michal Drozdal¹, Jakob Verbeek¹, Adriana Romero-Soriano^{1,2,4,5}

¹FAIR at Meta, ²Mila, Quebec AI Institute, ³Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK, France, ⁴McGill University, ⁵Canada CIFAR AI chair

*Core contributor

As the use of text-to-image generative models increases, so does the adoption of automatic benchmarking methods used in their evaluation. However, while metrics and datasets abound, there are few unified benchmarking libraries that provide a framework for performing evaluations across many datasets and metrics. Furthermore, the rapid introduction of increasingly robust benchmarking methods requires that evaluation libraries remain flexible to new datasets and metrics. Finally, there remains a gap in synthesizing evaluations in order to deliver actionable takeaways about model performance. To enable unified, flexible, and actionable evaluations, we introduce EvalGIM (pronounced “EvalGym”), a library for evaluating generative image models. EvalGIM contains broad support for datasets and metrics used to measure quality, diversity, and consistency of text-to-image generative models. In addition, EvalGIM is designed with flexibility for user customization as a top priority and contains a structure that allows plug-and-play additions of new datasets and metrics. To enable actionable evaluation insights, we introduce “Evaluation Exercises” that highlight takeaways for specific evaluation questions. The Evaluation Exercises contain easy-to-use and reproducible implementations of two state-of-the-art evaluation methods of text-to-image generative models: consistency-diversity-realism Pareto Fronts and disaggregated measurements of performance disparities across groups. EvalGIM also contains Evaluation Exercises that introduce two new analysis methods for text-to-image generative models: robustness analyses of model rankings and balanced evaluations across different prompt styles. In this paper, we outline the EvalGIM library and provide guidance for how others can add new datasets, metrics, and visualizations to customize the library for their own use cases. We also demonstrate the utility of EvalGIM by using its Evaluation Exercises to explore several research questions about text-to-image generative models, such as the role of re-captioning training data or the relationship between quality and diversity in early training stages. We encourage text-to-image model exploration with EvalGIM and invite contributions at <https://github.com/facebookresearch/EvalGIM/>.

Date: December 19, 2024

Correspondence: Melissa Hall at melissahall@meta.com

Code: <https://github.com/facebookresearch/EvalGIM/>

Blogpost: <https://ai.meta.com/blog/meta-fair-updates-agents-robustness-safety-architecture/>



1 Introduction

The rapid rise of text-to-image generative models has increased the practice of benchmarking their capabilities and weaknesses. Even as human evaluations have become an evaluation standard for text-to-image generative models, automatic evaluations remain common, as they allow for timely results about model performance, are easily scaled to many models and datasets, and can be reproduced given a standardized set-up. For example, it is standard to evaluate models with the Fréchet Inception Distance (FID) (Heusel et al., 2018) metric, and increasingly popular to benchmark with precision (Sajjadi et al., 2018b; Kynkäänniemi et al., 2019) and density (Naeem et al., 2020), which measure fidelity or “realness” of a set of generated images, and recall (Sajjadi et al., 2018b; Kynkäänniemi et al., 2019) and coverage (Naeem et al., 2020), which measure



EvalGIM

A Library for Evaluating Generative Image Models



Figure 1 EvalGIM (pronounced as "EvalGym") is an easy-to-use evaluation library for text-to-image generative models that unifies useful evaluation metrics, datasets, and visualizations, supports flexibility for user needs (and extensibility to future benchmarks), and provides actionable insights into model performance. To enable interpretable benchmarking, EvalGIM contains Evaluation Exercises that highlight takeaways for specific evaluation questions related to performance trade-offs, group representation, model ranking robustness, and prompting styles.

the diversity of generated images. Furthermore, there have been rapid advancements in benchmarks used for measuring image consistency (how well a generated image matches a text prompt), with alignment-based consistency metrics including CLIPScore (Hessel et al., 2021), VQAScore (Lin et al., 2024) and VIEScore (Ku et al., 2023) and question generation-based approaches, such as TIFA (Hu et al., 2023), Davidsonian Scene Graph (Cho et al., 2023) and VPEval (Cho et al., 2024).

While these metrics and benchmarks are useful, their lack of unification makes it inconvenient to perform robust, reproducible evaluations. In the natural language processing (NLP) domain, it has become quite common to perform multiple evaluations across different model skills with a unified benchmarking library (Wang, 2018; Wang et al., 2019; Liang et al., 2023). However, fewer multifaceted benchmarking libraries have been proposed for the text-to-image domain; The main benchmarking framework comprising multiple evaluation criteria is HEIM (Lee et al., 2023), which evaluates 12 different aspects, including image-text alignment and reasoning. Oftentimes, unified libraries present dozens of numbers which, while useful, require extensive study to identify interpretable takeaways about model development. This can make hypothesis-driven scientific studies challenging. Furthermore, they can become outdated quickly, as they rarely introduce newer metrics. With EvalGIM, we aim to address this gap by providing an easy-to-use evaluation library that unifies evaluation metrics, datasets, and visualizations, supports flexibility for user needs (and extensibility to future benchmarks), and provides actionable insights into model performance.

EvalGIM was designed to be used out-of-the-box, with immediate support for existing models and easy-to-run code targeted for specific evaluation tasks. To enable this, the library has a unified structure across image generation, evaluation, and visualization, so that resultant data is well-contained and easy to interact with. Furthermore, EvalGIM easily supports sweeps across different models or hyperparameters and evaluations disaggregated by subgroups of data with the use of a single flag, which provides a more thorough understanding of model performance. Given the large scale of generative image evaluations, EvalGIM also supports distributed evaluations across compute resources for faster analyses.

In addition, EvalGIM was designed for users to also contribute additional datasets and metrics to expand their analyses. We introduce inheritable `Dataset` types for real image data sources and prompt data sources to enforce consistent image, prompt, and metadata loading. These classes can be plugged directly into image generation and metric evaluation scripts, allowing for easy additions of new datasets regardless of their source. Similarly, EvalGIM builds on the `torchmetrics` library and defines a consistent set of wrapper functions that are used for updating and computing metrics and are compatible with the `Dataset` types. This makes it easy to add new metrics to the library such that they work out-of-the box.

Finally, EvalGIM contains "Evaluation Exercises" that allow for structured end-to-end evaluations targetting specific hypotheses about model performance. These Exercises can be executed in a reproducible manner

with user friendly code that requires only a list of model names to execute image generation and evaluation. The Exercises return publication-ready visualizations that usefully synthesize takeaways across datasets and metrics and clearly display key results that are relevant to the specific evaluation goals.

In this paper we start by outlining EvalGIM and its customizability for new metrics, datasets, and visualizations. Then, we introduce the Evaluation Exercises, which can be executed in a simple three-step process for ease and reproducibility. We apply each Evaluation Exercise to a preliminary research analysis, demonstrating how EvalGIM enables progress on text-to-image generative model research. For the first set of Evaluation Exercises, we reproduce evaluation methods of prior works:

Evaluation Exercise: Trade-offs. We study trade-offs in model performance via Pareto Fronts, as introduced by [Astolfi et al. \(2024\)](#). In our preliminary study with EvalGIM, we find that when training a text-to-image generative model, consistency increases steadily then plateaus, while automatic measures of quality and diversity can fluctuate.

Evaluation Exercise: Group Representation. We explore disaggregated representation measurements across geographic groups, as introduced by [Hall et al. \(2024b\)](#). In our early studies with EvalGIM, we find that advancements in latent diffusion models correspond to an improvement in quality and diversity more for some geographic regions (*e.g.*, Southeast Asia) than others (*e.g.*, Africa).

Additionally, EvalGIM contains Evaluation Exercises that present new analysis types:

Evaluation Exercise: Rankings Robustness. We study the robustness of model rankings across evaluation metrics and datasets. With EvalGIM, we highlight how the relative strengths of candidate models between quality *vs.* diversity can be obfuscated by the FID metric. In addition, consistency rankings across models can change depending on which metric is used, and model rankings are not always consistent across datasets.

Evaluation Exercise: Prompt Types. We perform analyses of balanced comparisons across datasets corresponding to different prompt styles. With EvalGIM, we find that mixing original and re-captioned training data can help improve diversity and consistency.

We hope that users find EvalGIM to be easy to use, with flexibility for custom needs and future benchmarks, and helpful in providing actionable insights into model performance. We release EvalGIM to enable future investigations of text-to-image generative models and encourage contributions of new datasets and metrics.

2 EvalGIM for Evaluating Generative Image Models

Throughout this Section, we describe the different text-to-image generative models, evaluation datasets, metrics, and visualizations supported in EvalGIM. We also describe how the library can easily be customized by the user to add new evaluation components or state-of-the-art benchmarks.

2.1 Models and sweeps

EvalGIM includes support for models that are accessible via the HuggingFace `diffusers` library. This includes functionality to support random seeding across all generations or custom seeding schemes, such as fixed seeds across prompts. Furthermore, the image generation, evaluation, and visualization scripts include support for hyperparameter sweeps, including over classifier guidance scales, allowing for a deeper understanding of nuances in model performance.

2.1.1 Add your own models

While any `diffusers` model can be run with EvalGIM out-of-the-box, the library can also be integrated directly into existing model training pipelines, allowing for more thorough monitoring of model performance over training time. This can be done by creating random latents to sample an existing text-to-image pipeline. See the repository `README` for pointers on how to adapt the code for new models.

2.2 Evaluation datasets

EvalGIM contains support for multiple evaluation datasets, including real image datasets and prompt-only datasets. We highlight currently supported datasets and how to customize the library by adding more datasets. For additional details about dataset construction and filtering, see Appendix A.

2.2.1 Real image datasets

Real image datasets have real-world images alongside human-written or automatically generated text prompts. EvalGIM includes standard text-to-image evaluation datasets **MS-COCO** (Lin et al., 2014), **ImageNet** (Deng et al., 2009), and **Conceptual 12M (CC12M)** (Changpinyo et al., 2021). Additionally, to provide insights into performance for different geographic regions, EvalGIM includes support for the **GeoDE** (Ramaswamy et al., 2023) dataset. GeoDE contains images of objects that were taken by people living in different countries around the world and has been used to evaluate text-to-image generative models for performance gaps across geographic regions (Hall et al., 2024b; Sureddy et al., 2024).

2.2.2 Prompt datasets

Prompt datasets contain only text used for conditioning image generation. EvalGIM supports the **PartiPrompts** (Yu et al., 2022) dataset containing 1600 prompts of varying complexity and theme. It also supports compositionality benchmarks **T2I-Compbench** (Huang et al., 2023b) and **DrawBench** (Saharia et al., 2022).

2.2.3 Add your own datasets

EvalGIM supports the addition of custom datasets. For additional real image datasets, developers can leverage the `RealImageDataset()` class, which requires a set of real images and optionally supports class- or group-level metadata. To incorporate new prompts used for image generation, developers may use the `RealAttributeDataset()` class, which requires prompt strings. This class also optionally supports class- or group-level labels and metadata corresponding to metric calculation, such as question-answer graphs.

```
1
2 class RealImageDataset(ABC):
3     """Dataset of real images, used for computing marginal metrics.
4     """
5
6     @abstractmethod
7     def __getitem__(self, idx) -> RealImageDatapoint:
8         """Returns RealImageDatapoint containing
9             image: Tensor
10            class_label: Optional[str]
11            group: Optional[List[str]]
12        """
13
14 class RealAttributeDataset(ABC):
15     """Dataset of prompts and metadata, used for generating images and computing metrics.
16     """
17
18     @abstractmethod
19     def __getitem__(self, idx) -> RealAttributeDatapoint:
20         """Returns RealAttributeDatapoint containing
21            prompt: str
22            condition: Condition
23            class_label: Optional[str]
24            group: Optional[List[str]]
25            dsg_questions: Optional[List[str]]
26            dsg_children: Optional[List[str]]
27            dsg_parents: Condition | None = None
28        """
```

Listing 1 New image and prompt datasets can easily be added to EvalGIM.

2.3 Metrics

EvalGIM provides support for many evaluation metrics. We highlight marginal metrics, conditional metrics, and grouped metrics included in EvalGIM and discuss how users can add their own metrics. We provide additional details about the metrics in Appendix B.

2.3.1 Marginal metrics: Image Realism & Diversity

Marginal metrics measure model performance by comparing distributions of images that are generated with prompts corresponding to real world images to the distribution of real images. We include **Fréchet Inception Distance** (FID) (Heusel et al., 2018), a standard evaluation metric that compares a distribution of generated images to a corresponding set of real images to provide an indicator of generative quality and diversity. To support more detailed insights into trade-offs in model performance, we include **precision** (Sajjadi et al., 2018a) and **density** (Kynkäänniemi et al., 2019) as indicators of generated image realism and **recall** (Sajjadi et al., 2018a) and **coverage** (Kynkäänniemi et al., 2019) as signals of generated image diversity.

2.3.2 Conditional metrics: Image-Text Consistency

Conditional metrics evaluate generated images while leveraging the prompt used in its conditioning, without any grounding with real-world images. EvalGIM includes consistency metrics that measure how well generated images match the text prompt used in its generation. **CLIPScore** (Hessel et al., 2021) embeds the generated image and text prompt and uses a CLIP model (Radford et al., 2021) to measure the cosine similarity between the two embeddings. **Davidsonian Scene Graph** (DSG) (Cho et al., 2023) leverages a language model to generate questions based on the text prompt and a visual question answering (VQA) model to answer the generated questions given the generated images. EvalGIM also includes the **VQAScore** (Lin et al., 2024), which uses a VQA model to predict the alignment between a text prompt and a generated image.

2.3.3 Grouped metrics

EvalGIM builds on prior works that study disparate performance between groups (Hall et al., 2024b; Sureddy et al., 2024) and automatically supports disaggregated group-level measurements. These evaluations are supported for all prompt datasets that contain group information, including subpopulations (*i.e.* geographic terms used when prompting), demographic groups (*i.e.* gender-presentation based on prompting that is used), and dataset partitions (*i.e.* color- or shape-specific prompts).

2.3.4 Add your own metrics

EvalGIM builds on `torchmetrics` (Detlefsen et al., 2022) as the framework for distributed metric calculations. Intermediate metric values are updated with each batch of real or generated images then computed holistically with the `compute()` function once all batches have been added. For marginal metrics, we introduce two primary functions for updating metric values at each batch: `update_real_images()` and `update_generated_images()` which take as inputs batches of real and generated images, respectively, and their associated metadata. Conditional metrics use only an `update()` function, which is equivalent to `update_generated_images()`.

```
1 class MarginalMetric(Metric):
2     def update_real_images(
3         self,
4         reference_images: torch.Tensor,
5         real_image_datapoint_batch: dict,
6     ) -> None:
7
8     def update_generated_images(
9         self,
10        generated_images: torch.Tensor,
11        real_attribute_datapoint_batch: dict
12    ) -> None:
13
14    def compute(self) -> dict:
15        return {"metric_name": super().compute()}
16
```

```

17 class ConditionalMetric(Metric):
18     def update(
19         self,
20         generated_images_batch: torch.Tensor,
21         real_attribute_datapoint_batch: dict
22     ) -> None:
23
24     def compute(self) -> dict:
25         return {"metric_name": super().compute()}

```

Listing 2 New metrics can easily be added to EvalGIM by leveraging the `torchmetrics` library.

Additionally, the EvalGIM supports the addition of metrics that can be disaggregated by subgroups. For information on how to adapt new metrics to support grouped measurements, see the README of EvalGIM.

2.4 Visualizations

The EvalGIM contains scripts for creating reproducible, easy-to-read visualizations.

2.4.1 Pareto Fronts

We implement Pareto Fronts visualizations, which were introduced in the context of understanding trade-offs in text-to-image model performance in [Astolfi et al. \(2024\)](#). These visualizations demonstrate the relationship between improvements of different metrics and can be used to easily plot many datapoints across model types or image generation hyperparameters. An example Pareto Front is shown in [Figure 3](#).

2.4.2 Radar plots

While Radar plots have been used for text-to-image model performance across different metrics ([BlackForest-Labs, 2024](#)), we extend their utility by leveraging them for group-level measurements. In EvalGIM, the radial axes correspond to group performance, and the relative location of groups along their axes allows for comparison of disparities in performance across groups. Furthermore, multiple models can be included on a plot to study which groups realize the most improvement. An example Radar plot is presented in [Figure 5](#).

2.4.3 Ranking table

We build on previous works ([Lee et al., 2023](#)) that present large tables of results across metrics and models by introducing ranking table visualizations. We also include multiple datasets for the same metric to provide insights into whether model rankings are consistent for different data distributions. EvalGIM provides additional interpretability by applying color coded indicators of model rankings across a given dataset and metric, visually revealing the robustness of a model’s ranking. An example is shown in [Figure 6](#).

2.4.4 Scatter plots

We also include scatter plots to compare performance across datasets that are subsampled to be the same size, enabling a fair comparison of marginal metrics between them. An example is presented in [Figure 7](#).

2.4.5 Add your own visualizations

Visualization scripts are stored in the `visualizations/` directory of EvalGIM and take as inputs a `csv` with evaluation results for each model-dataset-hyperparameter combination and possible visualization parameters, such as a list of metrics to display. New visualizations can be added following this scheme.

3 Analyzing Models with Evaluation Exercises

There are hundreds of different combinations of text-to-image model evaluations possible across dataset, metric, and visualization options in EvalGIM. A standard practice would be to either focus on a single dataset-metric-visualization combination for ease of actionability, at the cost of understanding the full picture

of model performance or generate a large table of metric results across different axes of evaluations, to the detriment of interpretability and actionability.

To address this gap, EvalGIM contains “Evaluation Exercises” for common analysis questions about text-to-image generative models. For example, a common evaluation question is, “What is the trade-off between quality, diversity, and consistency?” for a given model. Each Evaluation Exercise organizes a specific set of datasets, metrics, and visualizations to provide a clean, easy-to-interpret insight about text-to-image performance. While each Exercise could be extended with alternative datasets and metrics, we propose a specific subset here which yield reliable and interpretable analyses. To help with ease-of-use, each Evaluation Exercise runs with a simple notebook to generate images, evaluate, and visualize the main takeaways.

The Evaluation Exercises are summarized in Figure 2, and we describe them in further detail throughout this section. Additionally, we demonstrate their ease-of-use and interpretability by applying each one to a preliminary research analysis.





 Evaluation Exercise Goal	 Evaluation Dataset(s)	 Metric(s)	 Visualizations
#1 What is the trade-off b/w quality, diversity, & consistency?	COCO	Precision (Quality) Coverage (Diversity) VQAScore (Consistency)	Pareto Fronts
#2 Do models have strong group representation ?	GeoDE	Precision (Quality) Coverage (Diversity) CLIPScore (Consistency)	Radar plots
#3 Are model rankings robust ?	CC12M COCO GeoDE ImageNet	Precision (Quality) Coverage (Diversity) FID (Quality + Diversity) CLIPScore (Consistency) VQAScore (Consistency)	Ranking table
#4 For which prompt types do models perform best?	CC12M (15k samples) COCO (15k samples) GeoDE (15k samples) ImageNet (15k samples)	Precision (Quality) Coverage (Diversity) FID (Quality + Diversity) CLIPScore (Consistency)	Scatter plots

Figure 2 EvalGIM contains Evaluation Exercises that allow for structured end-to-end evaluations targeting specific analysis goals. The Exercises can be easily executed in a reproducible manner with user friendly notebooks.

3.1 Evaluation Exercise: Trade-offs

This Exercise explores tradeoffs across quality, diversity, and consistency of generated images. It is intended to validate and make reproducible analyses that have been previously explored by [Astolfi et al. \(2024\)](#).

3.1.1 Exercise definition

Datasets. Because this Evaluation Exercise leverages marginal metrics, it requires a dataset that contains real images. We use COCO, which is a standard in text-to-image generation evaluations and contains human-written, realistic captions that can be used as prompts.

Metrics. Alternatives to FID are used in order to differentiate between improvements in quality, diversity, and consistency. Precision is used to measure generated image quality, as it quantifies the proportion of generated images that fall within a manifold of real images. Coverage is used to quantify how well generated images represent the diversity of real images. To measure consistency, we depart from prior works and use VQAScore, which contains a higher correlation with human preference of image consistency compared to alternative metrics like CLIPScore ([Lin et al., 2024](#)).

Visualizations. Following [Astolfi et al. \(2024\)](#), we leverage Pareto Fronts drawn across different consistency, diversity, and quality metrics. These Pareto Fronts visualize trade-offs across different interventions, such as ablations (*e.g.* model size), hyperparameter sweeps (*e.g.* guidance scale), or training progress (*e.g.* iterations).

3.1.2 Exercise in action

As an example of the utility of the Trade-offs Evaluation Exercise, we use it to explore how a model varies in quality, diversity, and consistency throughout its early training process.

Experimental Set-up We train a text-to-image generative model that leverages flow matching (Lipman et al., 2023) and control conditions (Podell et al., 2023; Ifriqi et al., 2024) with a dataset of image-caption pairs including a subset from ImageNet (Deng et al., 2009), CC12M (Changpinyo et al., 2021), and an internally licensed dataset. The model is trained for a total of 650,000 iterations and evaluated with COCO using a classifier guidance scale of 7.5. Results are shown in Figure 3.

Findings We find that throughout the training process, quality, diversity, and consistency trade-off in different ways and are not always zero-sum. Consistency steadily increases initially in the training process, then plateaus at around 450,000 iterations. Quality (as measured by precision) can have a small but steady decrease, perhaps as generation capabilities increase beyond the distribution of the reference dataset used for evaluation. Furthermore, representation diversity shows an initial small increase then a slight decrease. Improvements in consistency coincide with small declines in quality. In addition, we inspect visual samples (see random examples in Figure 4) across the models and find that image-prompt alignment often visibly improves over the training process, aligning well with improvements in the consistency metric. Throughout training the images show improvements in finer-grained details, realistic shapes, and well-formed structures. It is left to future work to explore how these trends evolve over future iterations.

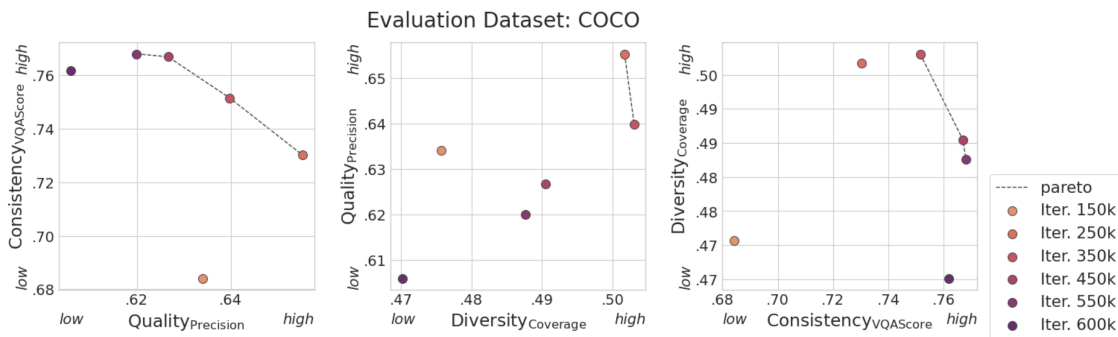


Figure 3 Utilizing the Trade-offs Evaluation Exercise gives insights into the relationship between quality, diversity, and consistency. When applied to preliminary studies of early training of a text-to-image generative model, consistency (as measured by VQAScore) increases then plateaus, while automatic measures of quality and diversity can fluctuate.



Figure 4 Finer-detailed image improvements identified via qualitative inspection do not always translate to improvements in automatic measures of quality and diversity.

3.2 Evaluation Exercise: Group Representation

This Evaluation Exercise studies performance of models for different subgroups of data, leveraging the disaggregated evaluation method presented by Hall et al. (2024b).

3.2.1 Exercise definition

Datasets We use the GeoDE dataset, disaggregated at the regional level. This corresponds to the geographic representation of different objects like “car” and “cooking pot” across regions around the world.

Metrics Precision and coverage are used to measure quality and diversity and CLIPScore for a measure of consistency. While prior work (Hall et al., 2024b) uses the bottom 10th-percentile CLIPScore, we focus on CLIPScore for the full prompt. We note that CLIPScore may give higher ratings to more stereotypical representations, as suggested in previous work (Agarwal et al., 2021; Hall et al., 2024a).

Visualization Radar plots are used for each metric, with axes corresponding to group performance. The relative location of groups along their axes allows for comparison of disparities in performance across group. Multiple models can be included on a plot to study which groups experience the most improvement, as indicated by the gap on a given group’s axis between model versions.

3.2.2 Exercise in action

Using the Group Representation Evaluation Exercise, we explore the evolution of geographic representation across different generations of the same family of latent diffusion models.

Experimental Set-up We evaluate different generations of a latent diffusion model (LDM). The first, LDM-1.5 (Rombach et al., 2022a), is trained on a public web dataset containing approximately 2 billion images and further trained on higher resolution images and fine-tuned on aesthetic images. The second, LDM-2.1 (Rombach et al., 2022a), is similarly trained on a dataset of approximately 5 billion images and further trained for multiple iterations on increasingly high-resolution samples, then fine-tuned. We also evaluate LDM-3 (Esser et al., 2024), which incorporates rectified flows (Liu et al., 2022; Albergo and Vanden-Eijnden, 2022; Lipman et al., 2023) and leverages three pre-trained text encoders. Finally, we evaluate LDM-XL (Podell et al., 2023), a base model to that generates latents using a larger UNet backbone and more parameters than the aforementioned models. We access all models via an API containing open-sourced model weights.

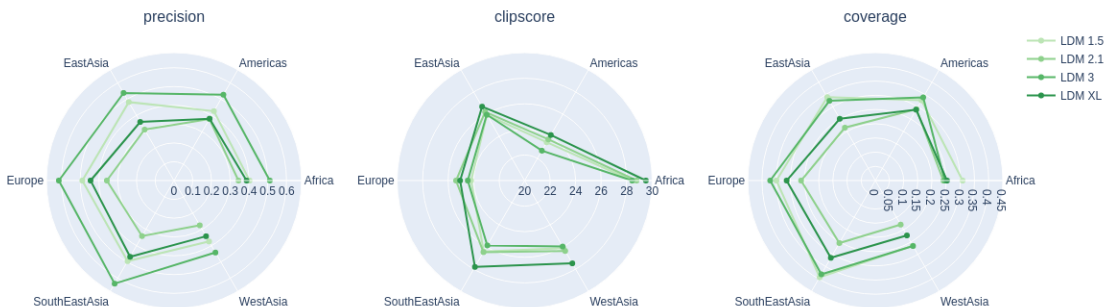


Figure 5 Using the Group Representation Evaluation Exercise provides insights into potential disparities in model performance across groups and whether improvements over successive model generations have occurred similarly across groups. When studying successive versions of a latent diffusion model with increasingly complex training data and fine-tuning methods, we find that advancements correspond to an improvement in quality and diversity more for some geographic regions (e.g. Southeast Asia) than others (e.g. Africa).

Findings With the Group Representation Exercise, we can study the gaps along each region’s radial axis to understand how modeling advances affect groups disparately. Figure 5 shows that there has been more improvement in the quality (precision) of images from Southeast Asia and Europe compared with those in West Asia and Africa. Notably, while previous work on earlier version of latent diffusion models highlighted *drops* in diversity and quality when studying geographic representation over successive versions of models (Hall et al., 2024b; Astolfi et al., 2024), we find the most recent model, LDM-3, *recovers* most of the performance loss of previous models. The one exception is for representational diversity for Africa, where LDM-3 lags behind the older LDM-1.5 in coverage measurements.

3.3 Evaluation Exercise: Ranking Robustness

We also introduce new evaluation formulations to improve understanding of text-to-image models. This Ranking Robustness Evaluation Exercise studies whether relative performance across different candidate models is consistent across a set of metrics and datasets. It is intended to serve as a more interpretable variant of large metric tables that are often provided in benchmarking efforts and proposals of new models, with a focus on highlighting whether claims of superior performance for a given model are actually consistent across different metrics and datasets.

3.3.1 Exercise definition

Datasets We recommend both well-established datasets and more recent variants that have been demonstrated to provide contextualization of standard benchmarks. For established datasets we include ImageNet and COCO. For newer datasets, we include GeoDE, which helps highlight whether overall improvements come at the cost of geographic representation, and CC12M, which uses alt-text captions as image prompts.

Metrics We recommend a combination of metrics that are included in other libraries, such as FID and CLIPScore (Lee et al., 2023), and additional, newer variants such as precision, coverage, and VQAScore to better contextualize performance claims. We include FID and additionally supplement with precision and coverage to demonstrate whether advances in quality or diversity, respectively, contribute to improvements in FID. For consistency, we include CLIPScore, which is more established, and the newer VQAScore, which has a stronger correlation with human annotations (Lin et al., 2024).

Visualizations This Evaluation Exercise is visualized with a ranking table across all datasets and metrics, where each cell is colored according to how the model ranks for the respective metric-dataset combination. Thus, it allows for easier identification of which models perform best across specific metrics or datasets and whether rankings across models are consistent.

3.3.2 Exercise in action

We study whether rankings of latent diffusion models are robust across different metrics and datasets.

Experimental set-up We study the same set of latent diffusion models described in Section 3.2: LDM-1.5, LDM-2.1, LDM-XL, and LDM-3.

	clipscore				fid_torchmetrics				coverage				precision				vqascore			
	CC12M	COCO	GeoDE	ImageNet	CC12M	COCO	GeoDE	ImageNet	CC12M	COCO	GeoDE	ImageNet	CC12M	COCO	GeoDE	ImageNet	CC12M	COCO	GeoDE	ImageNet
LDM XL	27.0	30.1	26.8	24.4	17.0	17.0	46.0	16.0	0.75	0.56	0.23	0.54	0.64	0.58	0.39	0.71	0.81	0.88	0.79	0.82
LDM 3	26.0	29.3	25.6	23.6	19.0	25.0	49.0	21.0	0.76	0.54	0.28	0.49	0.71	0.65	0.52	0.82	0.8	0.92	0.74	0.82
LDM 2.1	26.6	28.9	26.2	24.1	16.0	16.0	48.0	15.0	0.73	0.55	0.2	0.51	0.62	0.56	0.34	0.71	0.78	0.83	0.74	0.81
LDM 1.5	25.3	27.7	25.8	23.7	17.0	16.0	41.0	18.0	0.77	0.58	0.28	0.55	0.67	0.59	0.43	0.73	0.75	0.77	0.71	0.8

Figure 6 Utilizing the Ranking Robustness Evaluation Exercise across a suite of metrics and datasets reveals information that may be obfuscated by individual metrics. When applied to preliminary studies of latent diffusion models, we find that which models have better quality *v.s.* diversity can be obfuscated by FID. In addition, which model has the worst consistency can change depending on which consistency metric is used, and model rankings are not always consistent across datasets.

Findings The Rankings Robustness Evaluation Exercise provides insight into how information about model performance may be obfuscated across different metrics. We find that disaggregating measurements between quality (*i.e.* precision) and diversity (*i.e.* coverage) can tell a different story than combined measurements (*i.e.* FID). For example, we observe that LDM-3 ranks worst on FID across all datasets but highest for quality. Similarly, LDM-1.5 scores quite well on FID but this is likely due to its stronger diversity rather than quality. Furthermore, different consistency metrics can yield different outcomes of “best” models: when evaluating with GeoDE and ImageNet, LDM-1.5 shows the worst consistency when evaluated with VQAScore

while LDM-3 is worst when using CLIPScore (although the two are quite close). In addition, this analysis provides insights into the robustness of metrics across different datasets. For example, quality (*i.e.* precision) and consistency (*i.e.* CLIPScore and VQAScore) rankings correlate well across datasets. However, diversity (coverage) rankings for the models vary when evaluating with different datasets. For example, LDM-3 has among the strongest diversity when evaluated on the geographically representative GeoDE dataset but lowest on COCO and ImageNet.

3.4 Evaluation Exercise: Prompt Types

The Prompt Types Evaluation Exercise focuses on model performance across different prompt types, from a single concept word such as “apple” or “dog” to multi-sentence descriptions that capture details such as image styles, multiple objects, and layout. Analyzing text-to-image models along these different prompt types provides insights into whether a given intervention helps certain kinds of model interaction types (via prompting) more so than others.

3.4.1 Exercise definition

Datasets In a departure from prior works, we balance evaluation datasets so that they can be compared to each other with distributional metrics. In our case, we uniformly randomly subsample ImageNet, GeoDE, COCO, and CC12M so that they are the largest possible shared sized, e.g. approximately 15,000 images (as dictated by the size of the smallest dataset, CC12M). ImageNet prompts correspond to individual concepts, GeoDE prompts outline objects in a given geographic region, and COCO prompts are human-written captions of the corresponding real image. For CC12M, we use two types of prompts: the original alt-text based captions for the corresponding image and a re-captioned version using the outputs of the Florence-2 model (Xiao et al., 2023) when provided with the respective real image.

Metrics We use precision and coverage to measure quality and diversity and CLIPScore for consistency.

Visualization We include a scatterplot of precision *vs.* coverage across the evaluation datasets, which are comparable due to being balanced. We also include scatterplots of FID and CLIPScore.

3.4.2 Exercise in action

We use the Prompt Types Evaluation Exercise to perform an initial study of how different training data re-captioning methods can affect the diversity of generated images.

Experimental Set-up We focus on a latent diffusion model that has been trained with image-caption pairs from CC12M (Changpinyo et al., 2021) while employing different re-captioning strategies: Images are recaptioned with either PaliGemma (Beyer et al., 2024) (fine-tuned on COCO) or Florence-2 (Xiao et al., 2023).

Findings Overall, we observe that training models with both original and re-captioned data leads to higher precision and coverage performance across datasets. Furthermore, we find that re-captioning improves performance for some types of generation tasks more than others, with the stronger performance improvements occurring for the ImageNet and re-captioned CC12M datasets. In addition, we find that using only original captions is better than using only recaptioned data when generating images with the original caption distribution. Excluding original captions can have a larger negative effect on GeoDE and ImageNet. Furthermore, Florence2 recaptioning alone is worse than PaliCOCO recaptioning. When inspecting visual examples, we find that re-captioning training data with more dense images and using more dense captions at generation time can help with image consistency, such as increasing the presence of animals or people mentioned in the caption. However, using only re-captioned training data for a model can yield undesirable artifacts, like extraneous text or banners, when not prompting with similar prompt styles.

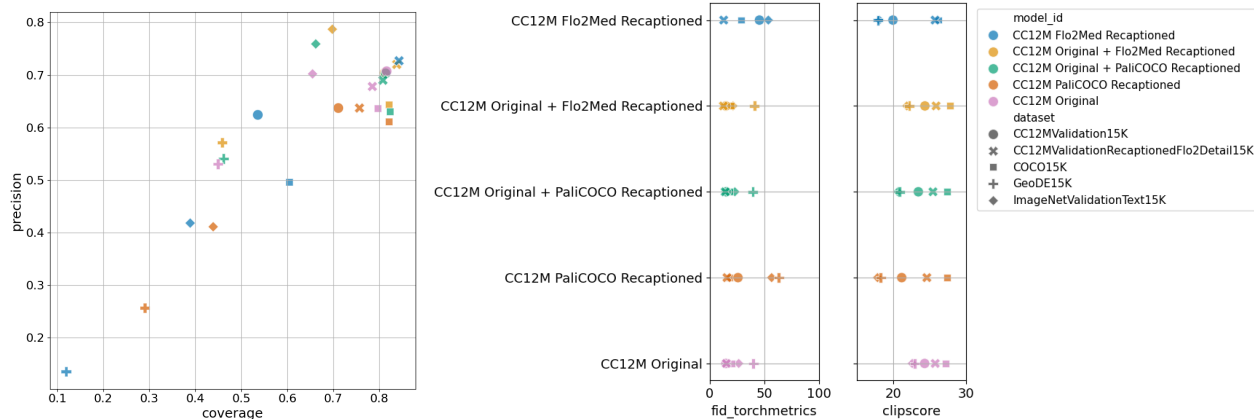


Figure 7 Analyzing models with the Prompt Types Evaluation Exercise provides insights into how training time interventions affect certain datasets more than others. When applied to a study of training text-to-image models with recaptured training data, we find that mixing original and recaptured training data can help improve diversity and consistency.



Figure 8 Generations with original CC12M captions (LEFT) and Flo2-Detailed recaptions (RIGHT). Using more descriptive captions when prompting the model can help increase the consistency of the images, even for models trained with coarser captions (e.g. CC12M Original model).

4 Limitations

While we encourage the use of EvalGIM to benchmark performance, we also discuss useful considerations related to automatic evaluation of text-to-image generative models.

Known limitations with automatic metrics Recent works have identified limitations of automatic metrics for text-to-image generative models, such as biases in backbone feature extractors that may prioritize textures (Geirhos et al., 2018), certain kinds of concepts (DeVries et al., 2019), or stereotypical representations of certain groups (Agarwal et al., 2021). In addition, a subset of work is focused on exploring the robustness of consistency metrics (Ross et al., 2024; Saxon et al., 2024). While automatic evaluations of text-to-image generative models are imperfect, they are useful in performing timely and large-scale benchmarking of

generative models. To enable increasingly robust automatic evaluations, the EvalGIM is designed to enable users to add new metrics as they are developed.

Dataset coverage Evaluation datasets are only as effective as their coverage of real world use-cases and people. Additional datasets may be added to the EvalGIM to increase evaluation coverage. These may build on benchmarks designed specifically for generative modeling to probe skills such as compositional understanding (Zhu et al., 2023; Huang et al., 2023a; Wu et al., 2024), commonsense reasoning (Fu et al., 2024) and disambiguation (Rassin et al., 2022). Furthermore, additional datasets that contain per-prompt group information, such as the FACET (Gustafson et al., 2023) dataset depicting people in different activities or professions, may be easily added to unlock new group measurements.

Incorporating human evaluations EvalGIM is not intended to replace human evaluations of text-to-image generative models or in-depth qualitative user studies. There are several existing frameworks designed for evaluations from humans, where users compare models head-to-head. Approaches such as GenAI Arena (Jiang et al., 2024), K-Sort Arena (Li et al., 2024b) and Text to Image Arena are dynamic and include humans-in-the-loop. GenAI-Bench (Li et al., 2024a) provides a large-scale user study using Likert scale to rate generated images along with correlations between humans and automatic metrics. Frameworks involving humans approaches are generally more fine-grained and reliable than automatic metrics, with the tradeoff of being more time consuming and costly.

5 Conclusion

In this work, we introduced EvalGIM, a library for evaluating generative models. EvalGIM complements existing libraries by including more recent metrics, reproducible visualizations, and Evaluation Exercises that focus benchmarking on specific analysis questions. In addition, the library has a strong focus on supporting customization with updated or new datasets, metrics, and visualizations to allow for adaptability to future state-of-the-art evaluation methods. In this paper, we leveraged EvalGIM to reproduce existing analysis methods focused on Pareto Front-based measurements of trade-offs and disaggregated measurements across geographic groups. Furthermore, we also introduced two new Evaluation Exercises focused on evaluating the robustness of model rankings and the effect of model interventions on different prompt styles. We hope that with the release of EvalGIM, researchers and practitioners will be able to reliably benchmark text-to-image models with actionable takeaways to guide future areas of development.

6 Acknowledgments

We thank Carolyn Krol for extensive consultation and support throughout this project. In addition, we also thank Brian Karrer and Matthew Muckley for their feedback and support throughout the project.

References

- Sandhini Agarwal, Gretchen Krueger, Jack Clark, Alec Radford, Jong Wook Kim, and Miles Brundage. Evaluating clip: Towards characterization of broader capabilities and downstream implications, 2021. <https://arxiv.org/abs/2108.02818>.
- Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. *arXiv preprint arXiv:2209.15571*, 2022.
- Pietro Astolfi, Marlene Careil, Melissa Hall, Oscar Mañas, Matthew Muckley, Jakob Verbeek, Adriana Romero Soriano, and Michal Drozdal. Consistency-diversity-realism pareto fronts of conditional image generative models, 2024. <https://arxiv.org/abs/2406.10429>.
- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023.
- Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, Thomas Unterthiner, Daniel Keysers, Skanda Koppula, Fangyu Liu, Adam Grycner, Alexey Gritsenko, Neil Houlsby, Manoj Kumar, Keran Rong, Julian Eisenschlos, Rishabh Kabra, Matthias Bauer, Matko Bošnjak, Xi Chen, Matthias Minderer, Paul Voigtlaender, Ioana Bica, Ivana Balazevic, Joan Puigcerver, Pinelopi Papalampidi, Olivier Henaff, Xi Xiong, Radu Soricut, Jeremiah Harmsen, and Xiaohua Zhai. Paligemma: A versatile 3b vlm for transfer, 2024. <https://arxiv.org/abs/2407.07726>.
- BlackForestLabs. Announcing black forest labs. <https://blackforestlabs.ai/announcing-black-forest-labs/>, 2024. [Online; accessed 10-December-2024].
- Andrew Brock. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, 2021.
- Jaemin Cho, Yushi Hu, Roopal Garg, Peter Anderson, Ranjay Krishna, Jason Baldridge, Mohit Bansal, Jordi Pont-Tuset, and Su Wang. Davidsonian scene graph: Improving reliability in fine-grained evaluation for text-image generation. *arXiv preprint arXiv:2310.18235*, 2023.
- Jaemin Cho, Abhay Zala, and Mohit Bansal. Visual programming for step-by-step text-to-image generation and evaluation. *Advances in Neural Information Processing Systems*, 36, 2024.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- Nicki Skaftø Detlefsen, Jiri Borovec, Justus Schock, Ananya Harsh Jha, Teddy Koker, Luca Di Liello, Daniel Stancu, Changsheng Quan, Maxim Grechkin, and William Falcon. Torchmetrics-measuring reproducibility in pytorch. *Journal of Open Source Software*, 7(70):4101, 2022.
- Terrance DeVries, Ishan Misra, Changhan Wang, and Laurens van der Maaten. Does object recognition work for everyone?, 2019. <https://arxiv.org/abs/1906.02659>.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis, 2024. <https://arxiv.org/abs/2403.03206>.
- Xingyu Fu, Muyu He, Yujie Lu, William Yang Wang, and Dan Roth. Commonsense-t2i challenge: Can text-to-image generation models understand commonsense? *arXiv preprint arXiv:2406.07546*, 2024.
- Shanghua Gao, Pan Zhou, Ming-Ming Cheng, and Shuicheng Yan. Masked diffusion transformer is a strong image synthesizer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23164–23173, 2023.
- Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *CoRR*, abs/1811.12231, 2018. <http://arxiv.org/abs/1811.12231>.

- Laura Gustafson, Chloe Rolland, Nikhila Ravi, Quentin Duval, Aaron Adcock, Cheng-Yang Fu, Melissa Hall, and Candace Ross. Facet: Fairness in computer vision evaluation benchmark. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 20370–20382, October 2023.
- Melissa Hall, Samuel J. Bell, Candace Ross, Adina Williams, Michal Drozdal, and Adriana Romero Soriano. Towards geographic inclusion in the evaluation of text-to-image models. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency, FAccT '24*, page 585–601, New York, NY, USA, 2024a. Association for Computing Machinery. ISBN 9798400704505. doi: 10.1145/3630106.3658927. <https://doi.org/10.1145/3630106.3658927>.
- Melissa Hall, Candace Ross, Adina Williams, Nicolas Carion, Michal Drozdal, and Adriana Romero Soriano. Dig in: Evaluating disparities in image generations with indicators for geographic diversity, 2024b. <https://arxiv.org/abs/2308.06198>.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium, 2018. <https://arxiv.org/abs/1706.08500>.
- Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20406–20417, 2023.
- Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A comprehensive benchmark for open-world compositional text-to-image generation. *Advances in Neural Information Processing Systems*, 36: 78723–78747, 2023a.
- Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A comprehensive benchmark for open-world compositional text-to-image generation, 2023b. <https://arxiv.org/abs/2307.06350>.
- Tariq Berrada Ifriqi, Pietro Astolfi, Melissa Hall, Reyhane Askari-Hemmat, Yohann Benchetrit, Marton Havasi, Matthew Muckley, Karteek Alahari, Adriana Romero-Soriano, Jakob Verbeek, and Michal Drozdal. On improved conditioning mechanisms and pre-training strategies for diffusion models, 2024. <https://arxiv.org/abs/2411.03177>.
- Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021. <https://doi.org/10.5281/zenodo.5143773>. If you use this software, please cite it as below.
- Dongfu Jiang, Max Ku, Tianle Li, Yuansheng Ni, Shizhuo Sun, Rongqi Fan, and Wenhu Chen. Genai arena: An open evaluation platform for generative models. *arXiv preprint arXiv:2406.04485*, 2024.
- Max Ku, Dongfu Jiang, Cong Wei, Xiang Yue, and Wenhu Chen. Viescore: Towards explainable metrics for conditional image synthesis evaluation. *arXiv preprint arXiv:2312.14867*, 2023.
- Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models. *Advances in neural information processing systems*, 32, 2019.
- Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models, 2019.
- Tony Lee, Michihiro Yasunaga, Chenlin Meng, Yifan Mai, Joon Sung Park, Agrim Gupta, Yunzhi Zhang, Deepak Narayanan, Hannah Benita Teufel, Marco Bellagente, Minguk Kang, Taesung Park, Jure Leskovec, Jun-Yan Zhu, Li Fei-Fei, Jiajun Wu, Stefano Ermon, and Percy Liang. Holistic evaluation of text-to-image models, 2023. <https://arxiv.org/abs/2311.04287>.
- Baiqi Li, Zhiqiu Lin, Deepak Pathak, Jiayao Li, Yixin Fei, Kewen Wu, Tiffany Ling, Xide Xia, Pengchuan Zhang, Graham Neubig, et al. Genai-bench: Evaluating and improving compositional text-to-visual generation. *arXiv preprint arXiv:2406.13743*, 2024a.
- Zhikai Li, Xuewen Liu, Dongrong Fu, Jianquan Li, Qingyi Gu, Kurt Keutzer, and Zhen Dong. K-sort arena: Efficient and reliable benchmarking for generative models via k-wise human preferences. *arXiv preprint arXiv:2408.14468*, 2024b.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Alexander Cosgrove, Christopher D Manning, Christopher Re, Diana Acosta-Navas, Drew Arad Hudson, Eric Zelikman,

- Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue WANG, Keshav Santhanam, Laurel Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Andrew Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. Holistic evaluation of language models. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. <https://openreview.net/forum?id=iO4LZibEqW>. Featured Certification, Expert Certification.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 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, 2014.
- Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation. *arXiv preprint arXiv:2404.01291*, 2024.
- Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling, 2023. <https://arxiv.org/abs/2210.02747>.
- Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- Muhammad Ferjad Naeem, Seong Joon Oh, Youngjung Uh, Yunjey Choi, and Jaejun Yoo. Reliable fidelity and diversity metrics for generative models. In *International Conference on Machine Learning*, pages 7176–7185. PMLR, 2020.
- William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
- Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis, 2023. <https://arxiv.org/abs/2307.01952>.
- Alec Radford, Jong Wook Kim, Chris Hallacy, A. Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- Vikram V. Ramaswamy, Sing Yu Lin, Dora Zhao, Aaron B. Adcock, Laurens van der Maaten, Deepti Ghadiyaram, and Olga Russakovsky. Geode: a geographically diverse evaluation dataset for object recognition, 2023. <https://arxiv.org/abs/2301.02560>.
- Royi Rassin, Shauli Ravfogel, and Yoav Goldberg. Dalle-2 is seeing double: flaws in word-to-concept mapping in text2image models. *arXiv preprint arXiv:2210.10606*, 2022.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, June 2022a.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022b.
- Candace Ross, Melissa Hall, Adriana Romero-Soriano, and Adina Williams. What makes a good metric? evaluating automatic metrics for text-to-image consistency. In *First Conference on Language Modeling*, 2024.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding, 2022. <https://arxiv.org/abs/2205.11487>.
- Mehdi S. M. Sajjadi, Olivier Bachem, Mario Lučić, Olivier Bousquet, and Sylvain Gelly. Assessing Generative Models via Precision and Recall. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018a.
- Mehdi SM Sajjadi, Olivier Bachem, Mario Lucic, Olivier Bousquet, and Sylvain Gelly. Assessing generative models via precision and recall. *Advances in neural information processing systems*, 31, 2018b.
- Michael Saxon, Fatima Jahara, Mahsa Khoshnoodi, Yujie Lu, Aditya Sharma, and William Yang Wang. Who evaluates the evaluations? objectively scoring text-to-image prompt coherence metrics with t2iscorescore (ts2). *arXiv preprint arXiv:2404.04251*, 2024.

- Abhishek Sureddy, Dishant Padalia, Nandhinee Periyakaruppa, Oindrila Saha, Adina Williams, Adriana Romero-Soriano, Megan Richards, Polina Kirichenko, and Melissa Hall. Decomposed evaluations of geographic disparities in text-to-image models, 2024. <https://arxiv.org/abs/2406.11988>.
- C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016.
- Alex Wang. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*, 2018.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems*, 32, 2019.
- Xindi Wu, Dingli Yu, Yangsibo Huang, Olga Russakovsky, and Sanjeev Arora. Conceptmix: A compositional image generation benchmark with controllable difficulty. *arXiv preprint arXiv:2408.14339*, 2024.
- Bin Xiao, Haiping Wu, Weijian Xu, Xiyang Dai, Houdong Hu, Yumao Lu, Michael Zeng, Ce Liu, and Lu Yuan. Florence-2: Advancing a unified representation for a variety of vision tasks, 2023. <https://arxiv.org/abs/2311.06242>.
- Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation, 2022. <https://arxiv.org/abs/2206.10789>.
- Xiangru Zhu, Penglei Sun, Chengyu Wang, Jingping Liu, Zhixu Li, Yanghua Xiao, and Jun Huang. A contrastive compositional benchmark for text-to-image synthesis: A study with unified text-to-image fidelity metrics. *arXiv preprint arXiv:2312.02338*, 2023.

Appendix

A Additional details: Datasets

We provide additional details about the datasets included in EvalGIM.

A.1 Real Image Datasets

We use the term *real image datasets* to define datasets that have real-world images alongside human-written or automatically generated text prompts. These datasets can be used for marginal or reference-free metrics.

COCO MS-COCO (Lin et al., 2014) is an image captioning dataset that is now used as a standard for text-to-image generative models (Betker et al., 2023). The EvalGIM supports COCO-2014 validation set, which contains approximately 40,000 images with 5 human-written captions per image. We default to using the first caption as the prompt for each generated image.

ImageNet ImageNet (Deng et al., 2009) is a standard classification dataset for computer vision tasks. It is useful for testing concept generation and utility of generations in downstream learning tasks (Gao et al., 2023; Peebles and Xie, 2023; Rombach et al., 2022b; Brock, 2018). The EvalGIM contains support for ImageNet for both class-conditional image generation models and text-to-image models (where the prompt corresponds to the image label).

CC12M The EvalGIM also contains the Conceptual 12M (CC12M) dataset (Changpinyo et al., 2021), which contains alt-text captions with up to 256 words per images.

GeoDE To provide insights into performance for different geographic regions, the EvalGIM includes support for the GeoDE (Ramaswamy et al., 2023) dataset. GeoDE contains images of objects that were taken by people living in different countries around the world and has been used to evaluate performance gaps across geographic regions (Hall et al., 2024b; Sureddy et al., 2024). We use the balanced version (Hall et al., 2024b) to ensure a “fair” comparison across objects and regions, since different quantities of images correspond to different manifold sizes and can impact measurements. This requires omitting objects that have too few samples and extraneous images of objects that had many samples and yields 29160 images, 180 images for each of the six regions for 27 objects. In addition, we remove images that contain the tree tag to avoid any spurious correlations with trees in our measurements.

A.2 Prompt Datasets

We define *prompt datasets* as datasets that are text-only. These can be used for reference-free metrics only.

PartiPrompts The PartiPrompts (Yu et al., 2022) dataset contains 1600 prompts of varying complexity and theme, spanning topics including “Imagination,” “Fine-grained detail,” “Writing and symbols,” and “Quantity.”

T2I-Compbench T2I-Compbench (Huang et al., 2023b) is a dataset of 3,000 prompts used to test compositionality skills of text-to-image generative models. It includes prompts corresponding to color-, shape-, and texture-binding, spatial- and non-spatial relationships, and complex compositions.

DrawBench The DrawBench (Saharia et al., 2022) dataset contains 200 prompts used to evaluate text-to-image models in categories including color, counting, (mis)spelling, rarity, positioning, and depictions of text.

B Additional details: Metrics

The library provides support for multiple evaluation metrics, including marginal ones, reference-free ones, and grouped ones, which we detail below.

B.1 Marginal metrics: Image Realism & Diversity

Marginal realism and diversity metrics measure model performance by generating images from prompts that correspond to real world images and compare the *marginal* distribution of generated images to the distribution of real images, as opposed to relying on image-to-image metrics.

FID: We include Fréchet Inception Distance (FID) (Heusel et al., 2018), a standard evaluation metric that compares the distribution of generated images to a corresponding set of real images. Frequently, the generated images are created with prompts that correspond to the real image dataset. FID is typically interpreted as an indicator of both image quality and diversity.

Precision/Coverage/Recall/Density: To support more detailed insights into trade-offs in model performance, we include precision (Sajjadi et al., 2018a) and density (Kynkäänniemi et al., 2019) metrics to provide an indicator of generated image realism. Precision quantifies the proportion of generated images that fall within the manifold of real images while density additionally accounts for the *quantity* of real images the generated images are close to. We include recall (Sajjadi et al., 2018a) and coverage (Kynkäänniemi et al., 2019) as indicators of the diversity of generated images. Recall corresponds to the proportion of real images that fall in the manifold of generated images, while coverage measures the proportion of real images whose manifolds contain at least one generated image. The EvalGIM includes support for the InceptionV3 (Szegedy et al., 2016) model as the feature extractor for computing these metrics and uses a nearest-neighbors value of $k = 3$ to construct manifolds.

B.2 Conditional metrics: Image-Text Consistency

Reference-free metrics evaluate generated images without any grounding on real-world images. The EvalGIM includes consistency metrics that compute a score to measure how well a (generated) image I matches a given text prompt T . These metrics do not require a set of reference images.

CLIPScore: CLIPScore (Hessel et al., 2021) embeds the generated image and text prompt and uses a CLIP model (Radford et al., 2021) to measure the cosine similarity between the two embeddings. As a note, CLIPScore may “reward” more stereotypical representations (Hall et al., 2024a). The EvalGIM includes support for the OpenCLIP (Ilharco et al., 2021) model.

Davidsonian Scene Graph (DSG): DSG (Cho et al., 2023) leverages a language model to generate questions based on the text prompt and a visual question answering (VQA) model to answer the generated questions given the generated images.

VQAScore: VQAScore (Lin et al., 2024) uses a VQA model to predict the alignment between a text prompt and a generated image. It does so by computing the probability of the answer “yes” to the question, “Does this figure show {text}?”