

# Multilingual Conceptual Coverage in Text-to-Image Models

Michael Saxon

University of California, Santa Barbara saxon@ucsb.edu

William Yang Wang

University of California, Santa Barbara william@cs.ucsb.edu



Figure 1: A selection of images generated by DALLE-mega, Stable Diffusion 2, DALLE-2, and AltDiffusion, illustrating their *conceptual coverage* of “dog,” “airplane,” and “face” across English, Spanish, German, Chinese (simplified), Japanese, Hebrew, and Indonesian. Coverage of the concepts varies considerably across model and language, and can be observed in the consistency and correctness of images generated under simple prompts.

## Abstract

We propose “Conceptual Coverage Across Languages” (*CoCo-CroLa*), a technique for benchmarking the degree to which any generative text-to-image system provides multilingual parity to its training language in terms of tangible nouns. For each model we can assess “conceptual coverage” of a given target language relative to a source language by comparing the population of images generated for a series of tangible nouns in the source language to the population of images generated for each noun under translation in the target language. This technique allows us to estimate how well-suited a model is to a target language as well as identify model-specific weaknesses, spurious correlations, and biases without a-priori assumptions. We demonstrate how it can be used to benchmark T2I models in terms of multilinguality, and how despite its simplicity it is a good proxy for impressive generalization.

## 1 Introduction

Neural text-to-image models convert text prompts into images (Mansimov et al., 2016; Reed et al., 2016) using internal representations reflective of the training data population. Advancements in conditional language modeling (Lewis et al., 2020), variational autoencoders (Kingma and Welling, 2014), GANs (Goodfellow et al., 2020), multimodal representations (Radford et al., 2021), and latent diffusion models (Rombach et al., 2022) have given rise to sophisticated text-to-image (T2I) systems that exhibit impressive *semantic generalization capabilities*, with which they generate coherent, visually-appealing images with novel combinations of objects, scenarios, and styles (Ramesh et al., 2021). Their semantic latent spaces (Kwon et al., 2022) ground words to associated visuals (Hutchinson et al., 2022). Characterizing the limits of these systems’ capabilities is a challenge. They



Figure 2: We hypothesize that a model’s ability to generate creative, compositional images depicting tangible concepts (e.g., astronaut, horse, soup, bear) is predicated on its ability to generate simple images of the concepts alone. Samples from Ramesh et al. (2022).

are composed of elements trained on incomprehensibly large (Prabhu and Birhane, 2020; Jia et al., 2021), web-scale data (Gao et al., 2021; Schuhmann et al., 2021), hindering training-data-centric model analysis (Mitchell et al., 2019; Gebru et al., 2021) to address this problem.

Demonstrations of novel T2I model capabilities tend to rely on the subjective impressiveness of their ability to generalize to complex, novel prompts (Figure 2). Unfortunately, the space of creative prompts is in principle infinite. However, we observe that **impressive creative prompts are composed of known, tangible concepts**.

Can we directly evaluate a model’s knowledge of these tangible concepts as a partial proxy for its capability to generalize to creative novel prompts? Perhaps. But finding a diverse set of significant failure cases of basic concept knowledge for these models is challenging—in their training language.

We observe that when prompted with simple requests for specific tangible concepts in a constrained style, T2I models can sometimes generate consistent and semantically-correct images in languages for which they received limited training (Figure 1, Figure 3). We refer to this capacity as *concept possession* by said model in said language. At scale, we can assess the language-concept possession for a diverse array of concepts and languages in a model to attempt to describe its overall multilingual generalization capability. We refer to the degree of this capability as the model’s multilingual *conceptual coverage*. In this work we:

1. Introduce objective measures of *multilingual conceptual coverage* in T2I models that compare images generated from equivalent prompts under translation (Figure 4).
2. Release **CoCo-CroLa**, a benchmark set for conceptual coverage testing of 193 tangible concepts across English, Spanish, German, Chinese, Japanese, Hebrew, and Indonesian.

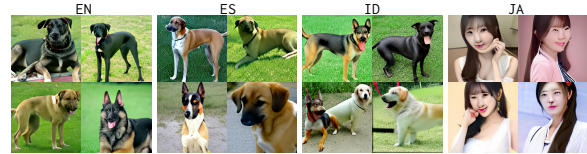


Figure 3: Although DALL-E mini (Dayma et al., 2021) is ostensibly trained only on English data, when elicited with “big dog” in Spanish, Indonesian, and Japanese it generalizes the “dog” concept to ES and ID, while exhibiting an offensive concept-level collision in JA.

3. Validate the utility of *conceptual coverage analysis* with a preliminary pilot study suggesting that generalization to complex, creative prompts follows concept possession.

Our benchmark enables fine-grained concept-level model analysis and identification of novel failure modes, and will guide future work in increasing the performance, explainability, and linguistic parity of text-to-image models.

## 2 Motivation & Related Work

This work is an attempt to produce a scalable, precise technique for characterizing conceptual coverage with **minimal assumptions** about the concepts or models themselves. In this section we lay out our motivations alongside relevant related work.

**Benchmarks enabling model comparability have been a driving force** in the development of pretrained language models (LM) (Devlin et al., 2019). For classification and regression tasks, evaluation under fine-tuning (Howard and Ruder, 2018; Karpukhin et al., 2020) is a straightforward and practical proxy for pretrained LM quality (Dodge et al., 2020) (e.g., for encoder-only transformer networks (Liu et al., 2019)). For these classification models, higher performance on benchmark datasets (Lai et al., 2017; Rajpurkar et al., 2018; Wang et al., 2019) became the epitome of LM advancement. However, other important qualities in models including degree of social biases (Sheng et al., 2019) and robustness (Clark et al., 2019) arising from biases in training data (Saxon et al., 2021) can only be captured by more sophisticated benchmarks that go beyond simple accuracy (Cho et al., 2021). CheckList represented an influential move in this direction by benchmarking model performance through *behavioral analysis under perturbed elicitation* (Ribeiro et al., 2020).

In contrast, generative large language models (LLMs) such as GPT-3 (Brown et al., 2020) have

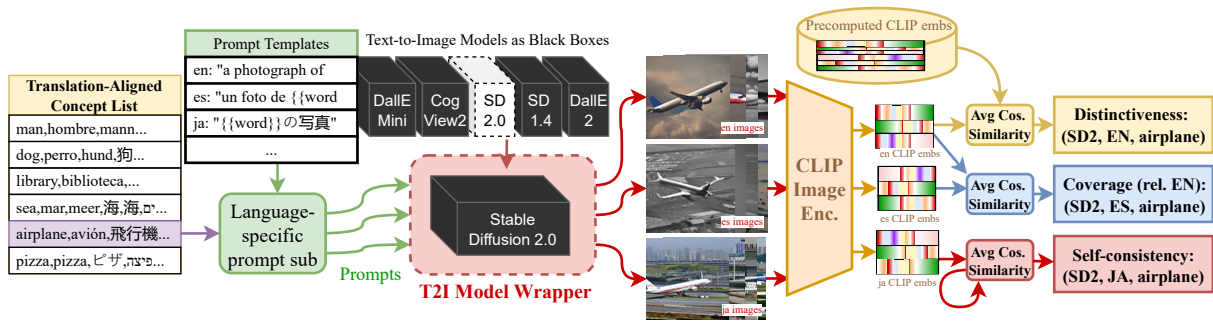


Figure 4: **CoCo-CroLa** assesses the cross-lingual coverage of a concept in a model by plugging all the term translations into prompt templates, generating a set of images from a model under test, extracting their corresponding CLIP embeddings, and computing concept-level **distinctiveness**, **coverage**, and **self-consistency** for the concept with respect to each language. (Demo and code available at [github.com/michaelsaxon/CoCoCroLa](https://github.com/michaelsaxon/CoCoCroLa))

a broader range of outputs, use-cases, and capabilities, making evaluation more difficult. For many text-generative tasks such as summarization and creative text generation, the crucial desired quality is subjective, and challenging to evaluate (Xu et al., 2022). However, as these LLMs operate in a text-only domain, existing supervised tasks could be ported to few-shot or zero-shot evaluations of LLM capabilities (Srivastava et al., 2022). While performance on these benchmarks isn’t directly indicative of the impressive generative performance and generalization capabilities, they are a means to measure improvement (Suzgun et al., 2022).

**Text-to-image models are even more difficult to evaluate than LLMs.** Unlike in LLMs, there are no objective evaluation tasks that can be directly ported as proxy evaluations. For example, while GPT-3 was introduced with impressive zero-shot performance across many classification tasks, the T2I model DALL-E 2 was primarily introduced with human opinion scores and cool demo images (Ramesh et al., 2022). Prior efforts in developing T2I evaluations such as Drawbench (Saharia et al., 2022) and DALL-Eval (Cho et al., 2022) fall into the trap of trying to build “everything benchmarks” for which whether the benchmark accurately reflects the practical task being asked of the computer in its real-world context is difficult to assess (Raji et al., 2021). We instead seek to build an *atomic benchmark* which narrowly and reliably captures a specific characteristic—conceptual knowledge as reflected by a model’s ability to reliably generate images of an object.

**Multilingual conceptual coverage is a high-variation T2I model performance setting.** (Figure 1) Perhaps more importantly, it has immediate

value, as work on improving T2I model multilinguality has been proposed, but hampered by a lack of evaluation metrics.

Chen et al. (2022) introduce AltCLIP and Alt-Diffusion, models produced by performing multilingual contrastive learning on a CLIP checkpoint for an array of non-English languages including Japanese, Chinese, and Korean. Without an objective evaluation benchmark, they can only demonstrate their improvement through human evaluation of impressive but arbitrary examples. **CoCo-CroLa** improves this state of affairs by enabling Checklist-like direct comparison of techniques for reducing *multilingual conceptual coverage* disparities as an objective, capabilities-based benchmark.

Excitingly, we find that **conceptual coverage is upstream of the impressive T2I model creativity** that model developers and end-users are fundamentally interested in. This means that not only is **CoCo-CroLa** an objective evaluation of T2I system capabilities, it is also a **proxy measure for the deeper semantic generalization capabilities we are interested in enhancing** in second languages, as we demonstrate in subsection 5.6.

### 3 Definitions & Formulations

We define a multilingual *concept* over languages  $L$  as a set of words in each language carrying the same meaning and analogous colloquial use. We refer to the equivalent translation of concept  $c_k$  in language  $\ell$  as  $c_{k,\ell}$ .

Given a set of concepts  $C$ , test language  $\ell$ , a *minimal eliciting prompt*<sup>1</sup>  $MP_\ell$ , text-to-image model

<sup>1</sup>We define a minimal eliciting prompt as a short sentence with a slot for concept work insertion, intended to enforce style consistency without interfering with the concept.

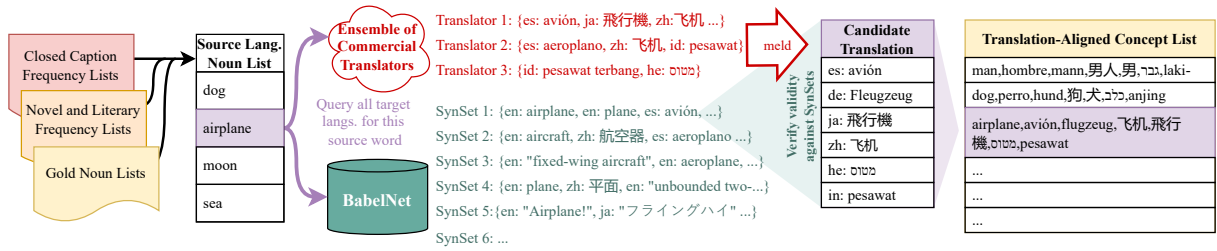


Figure 5: A diagram of our approach for producing the aligned noun concept list across the target language set using an ensemble of cloud translation services and BabelNet. Full description of this method in [Appendix A](#).

$f$ , and a desired number of images-per-concept  $n$ , we sample  $n|C||L|$  images  $I_{c_k,\ell,i}$ , where

$$I_{c_k,\ell,i} \sim f(\text{MP}_\ell(c_k,\ell)) \quad (1)$$

For every concept word in the language  $\ell$   $c_k \in C$ .

Given an image feature extractor  $F$ , some similarity function  $\text{SIM}(\cdot, \cdot)$ , we assess whether  $f$  possesses concept  $c_k,\ell$  in the test language if using the following metrics from the concept-image set  $\{I_{c_k,\ell,i}\}_{i=0}^n$  (**CoCo-CroLa** scores in [Figure 4](#)):

**Distinctiveness.** The images are *distinct* if they tend to not resemble the population of images generated for other concepts in the target language.

Formally, we compute our (inverse) distinctiveness score  $\text{Dt}(f, \ell, c_k)$  relative to  $m$  images sampled from other concepts in  $C$ :

$$\text{Dt} = \frac{1}{mn} \sum_{j=0}^m \sum_{i=0}^n \text{SIM}(F(I_{c_k,\ell,i}), F(I_{c_r,\ell,s})), \quad (2)$$

$$c_{r,\ell} \sim C \setminus c_{k,\ell}, \quad s \sim U\{0, n\} \quad (3)$$

**Self-consistency.** The images are *self-consistent* if they tend to resemble each other as a set.

Formally, we compute the self-consistency score  $\text{Sc}(f, \ell, c_k)$  as:

$$\text{Sc} = \frac{1}{n^2 - n} \sum_{j=0}^n (\sum_{i=0}^n \text{SIM}(F(I_{c_k,\ell,i}), F(I_{c_k,\ell,j})) - 1) \quad (4)$$

We subtract 1 from each step in the numerator and  $n$  from the denominator so that identical matches generated image to itself.

**Correctness.** The images are *correct* if they faithfully depict the object being described.

Rather than assess this using a classification model (hindering generality depending on the pretrained classifier), we use faithfulness relative to a source language  $\ell_s$ , cross consistency  $\text{Xc}(f, \ell, c_k, \ell_s)$  as a proxy:

$$\text{Xc} = \frac{1}{n^2} \sum_{j=0}^n \sum_{i=0}^n \text{SIM}(F(I_{c_k,\ell,i}), F(I_{c_k,\ell_s,j})) \quad (5)$$

It is important to note that in the source language (eg. English),  $\text{Xc} \approx \text{Sc}$ , as for that language both metrics are essentially comparing the same (concept, language) pair to itself.

Thus we augment with a second language-grounded correctness score, utilizing the average text-image similarity score of the English concept text against the set of generated images, for a CLIP image encoder  $F$  and text encoder  $F_t$ ,  $\text{Wc}$ :

$$\text{Wc} = \frac{1}{n} \sum_{i=0}^n F_t(c_{k,\ell_s}) \cdot F(I_{c_k,\ell,i}) \quad (6)$$

## 4 Approach

We compute distinctiveness, self-consistency, and correctness scores across English, Spanish, German, Chinese (Simplified), Japanese, Hebrew, and Indonesian on the models listed in [Table 1](#).

For each (language, concept) pair, we generate 10 images for analysis. We use a CLIP ([Radford et al., 2021](#)) checkpoint from HuggingFace<sup>2</sup> as our semantic visual feature extractor  $F$ , and cosine similarity as our similarity function ( $\text{SIM}(\mathbf{a}, \mathbf{b}) = \mathbf{a} \cdot \mathbf{b} / \|\mathbf{a}\| \|\mathbf{b}\|$ ). We collect a translation-aligned concept list  $C$  using techniques described in [subsection 4.1](#) and depicted in [Figure 5](#). We release our list generation code, testing code, feature extraction code, and final concept list as **CoCo-CroLa v1.0**<sup>3</sup>.

### 4.1 Translation-aligned concept set collection

We automatically produce an aligned multilingual concept list, where meaning, colloquial usage, and connotations are preserved as well as possible. We

<sup>2</sup>HF: [openai/clip-vit-base-patch32](https://openai.com/clip-vit-base-patch32).

<sup>3</sup>Demo and code: [github:michaelsaxon/CoCoCroLa](https://github.com/michaelsaxon/CoCoCroLa)

Model	Authors (Year)	Repository	Training Language
DALL-E Mini	Dayma et al. (2021)	github:borisdayma/dalle-mini	EN
DALL-E Mega			
CogView 2	Ding et al. (2021)	github:THUDM/CogView	ZH
StableDiffusion 1.1	Rombach et al. (2022)	HF:CompVis/stable-diffusion-v1-1	EN
StableDiffusion 1.2		HF:CompVis/stable-diffusion-v1-2	
StableDiffusion 1.4		HF:CompVis/stable-diffusion-v1-4	No language filter
StableDiffusion 2		HF:stabilityai/stable-diffusion-2	
DALL-E 2	Ramesh et al. (2022)	openai.com/dall-e-2/ (no checkpoints)	No language filter
AltDiffusion m9	Chen et al. (2022)	HF:BAAI/AltDiffusion-m9	EN, ES, FR, IT, RU, ZH, JA, KO

Table 1: The set of text-to-image models we evaluated with *CoCo-CroLa v1.0*. Some monolingual models may integrate pretrained elements such as CLIP checkpoints that have been trained on multilingual data.

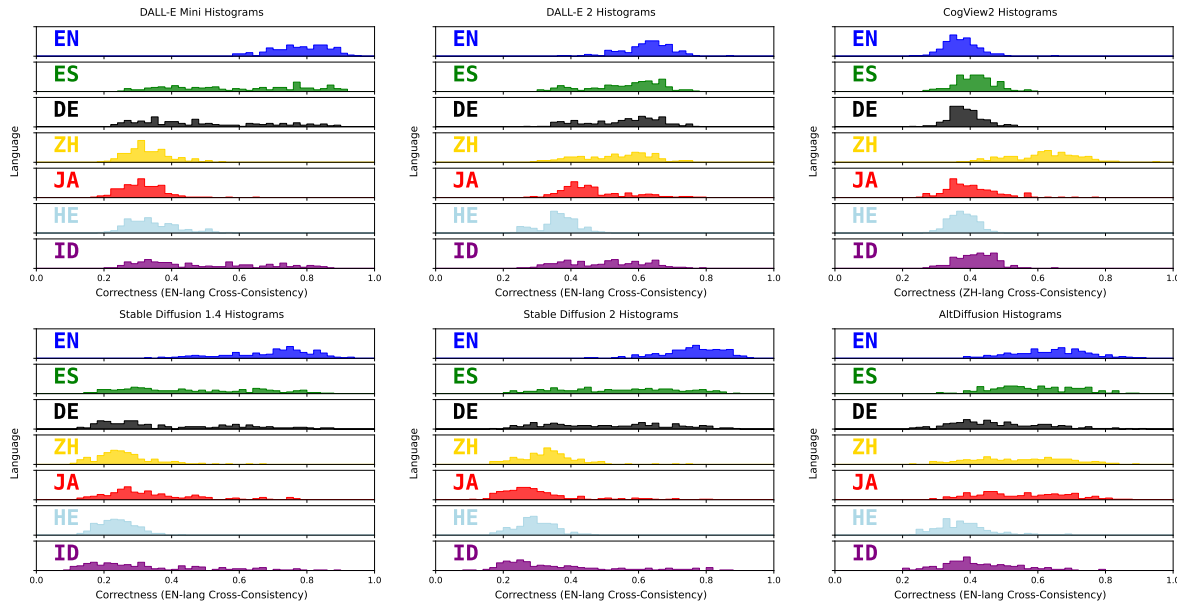


Figure 6: Histograms of the distribution of **correctness** cross-consistency ( $X_c$ ) for each test language for six assessed models. Rightward probability mass reflects better conceptual coverage.

identify *tangible nouns* describing physical objects, animals, people, and natural phenomena as a class of concepts that are both straightforward to evaluate and tend toward relative ubiquity in presence as words across languages and cultures, to the extent this ubiquity can be ensured through automated means.

Automated production is desirable for this task, as it enables new languages to be easily added in the future. To minimize translation errors we utilize both a large knowledge graph of terminology and an ensemble of commercial machine translation systems to produce an aligned concept list (Figure 5). Full details for our translation pipeline, as well as the full concept list, are in Appendix A.

## 4.2 Making minimal eliciting prompts

As discussed in section 3, an ideal prompt template would enforce stylistic consistency in the generated outputs without introducing biases that interfere with the demonstration of concept possession.

Following Bianchi et al. (2022) we build simple prompts of the form, “a photograph of \_\_\_\_\_”, which we manually translate into target languages. This simple template-filling approach will introduce grammatical errors for some languages. We briefly investigate if this matters in Appendix B.

## 4.3 Applying the metrics for analysis

We assess Dt, Sc, Xc, and Wc for each (concept, language) pair for each model. Using these we compare models and assess the validity of conceptual coverage as a proxy for generalization.

## 5 Findings

Figure 6 shows histograms for the distributions of the **cross-consistency correctness** proxy score  $X_c$  for each concept, relative to the training language (either English or Chinese) for DALL-E Mini, DALL-E 2, CogView 2, Stable Diffusion 1.4, Stable Diffusion 2, and AltDiffusion across the

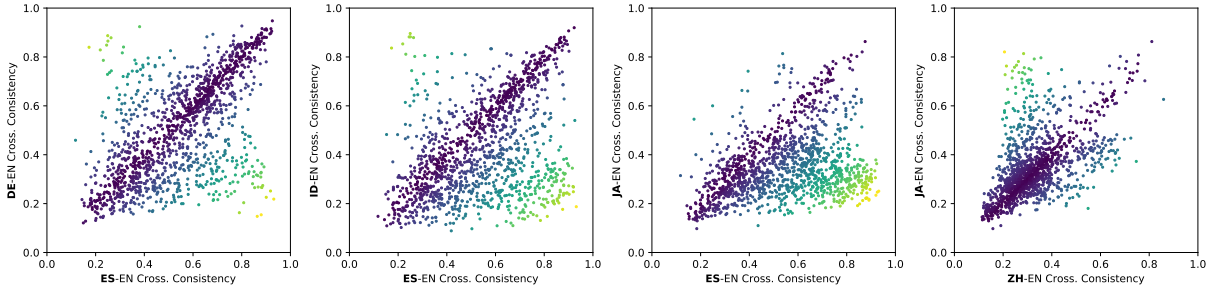


Figure 7: The **correctness** score for every (concept, model) pair for (right to left) ES vs DE, ES vs ID, ES vs JA, and ES vs JA. Languages sharing scripts (ES/DE/ID and JA/ZH) are more correlated than those that don't (ES/JA).

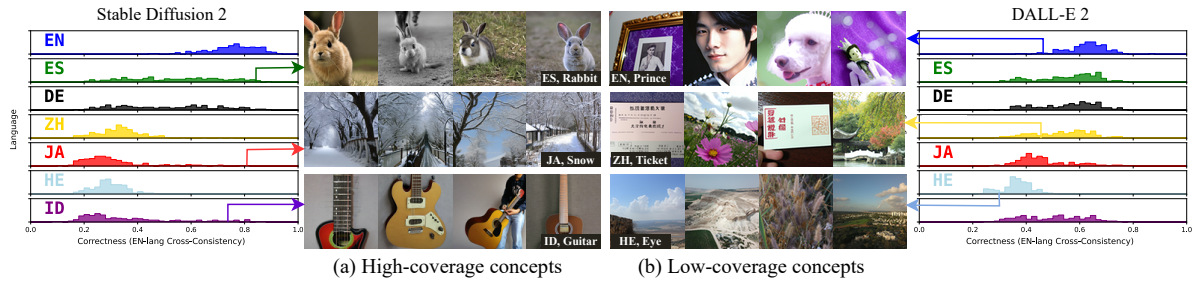


Figure 8: We automatically identify (a) high-coverage concepts in Stable Diffusion 2 (ES, rabbit), (JA, snow), (ID, guitar) and (b) low-coverage concepts in DALL-E 2 (EN, prince), (ZH, ticket), (HE, eye) using **correctness** Xc.

seven test languages. This plot clearly depicts that for the primarily English-trained models (DALL-E Mini, Stable Diffusion 1.4, Stable Diffusion 2), English-language performance (a high-mean distribution of high-EN-EN consistency concepts) is considerably better than the other languages. Similarly, for CogView2, trained on Chinese, the Chinese distribution of ZH-ZH scores is considerably better than the others, which do equally bad.

DALL-E 2 received open-ended multilingual training, and exhibits more consistent acceptable performance across the European and East Asian languages being tested. AltDiffusion, which has had its CLIP text encoder contrastively trained against multilingual representations on 9 languages (including ES, DE, ZH, and JA) exhibits higher performance on its training languages than its non-training languages (HE and ID).

Correctness distributions for Spanish, German, and Indonesian look roughly similar (in terms of mean and variance) for all models but AltDiffusion. This is particularly interesting because they are the three non-English languages that also use the Latin alphabet. Figure 7 compares the **correctness** Xc score for every concept, in every model, across pairs of languages that fully or partially share scripts (ES, DE, ID), (ZH, JA) and two languages that don't (JA, ES). Across pairs of lan-

guages that share scripts, there is a high correlation between possession of a given concept in one language and the other. A consistent trend across all models was poor performance on Hebrew, which is both considerably lower-resource compared to the other six test languages, and uses its own unique writing system.

### 5.1 Correctness feature captures possession

Figure 8 shows how choosing samples of an image generated by a model, elicited by a high- or low-**correctness** score naturally reveals in which languages which concepts are possessed (e.g., Stable Diffusion 2 possesses ES:rabbit, JA:snow, and ID:guitar.) When a model possesses a concept, the outputted images are often visually similar with the tangible concept set in similar scenarios.

### 5.2 Types of concept non-possession

A model *not* possessing a concept can manifest in a few different scenarios depicted in Figure 8 (b). DALL-E 2 doesn't possess "prince" in English because it outputs a variety of different images, including human portrait photos, and pictures of photos, toys, and dogs. The existence of these *non-specific* error cases reflects the imperfect nature of our automated concept collection pipeline. Removing these poorly specified concepts is one way we plan to improve *CoCo-CroLa* in future releases.

Model	EN		ES		DE		ZH		JA		HE		ID		Avg	
	Xc	Wc	Xc	Wc	Xc	Wc	Xc	Wc	Xc	Wc	Xc	Wc	Xc	Wc	Xc	Wc
DALL-E Mega	<b>81</b>	<b>28</b>	<b>65</b>	26	<b>64</b>	<b>26</b>	29	21	32	21	28	19	<b>51</b>	<b>25</b>	50	<b>24</b>
DALL-E Mini	78	27	59	25	50	23	33	21	31	21	34	<b>20</b>	49	24	48	23
SD 1.1	69	26	52	23	46	22	32	19	37	21	28	17	39	21	43	21
SD 1.2	71	26	48	23	44	22	28	19	35	21	24	17	37	21	41	21
SD 1.4	69	26	46	23	40	22	26	20	34	21	24	17	34	21	39	21
SD 2	76	27	54	24	51	24	34	19	31	21	29	17	37	21	45	22
CogView 2	37	20	42	20	39	20	<b>62</b>	<b>25</b>	40	21	<b>38</b>	<b>20</b>	42	20	43	21
DALL-E 2	61	27	55	<b>27</b>	54	<b>26</b>	44	<b>25</b>	42	22	36	19	42	<b>25</b>	48	24
<b>AltDiffusion m9</b>	64	26	59	25	49	22	55	<b>25</b>	<b>55</b>	<b>25</b>	<b>38</b>	<b>20</b>	43	22	<b>52</b>	23
Avg	67	26	53	24	49	23	38	22	38	21	31	18	42	22		

Table 2: **Correctness** scores (Xc and Wc) averaged for all concepts within a column language for all models. Note that Xc for CogView2 is relative to ZH rather than EN. AltDiffusion performs best in terms of total average Xc, and number of Xc or Wc column “wins.” DALL-E Mega performs best on Latin languages and average Wc.

A second type of possession failure we observe, we dub *specific collisions*. For example, Figure 1 and Figure 3 show JA collisions for the DALL-E mini/mega family. Both models consistently generate images of humans for “dog” but pictures of green landscape scenes for “airplane.” While these generated concepts are incorrect, they represent an incorrect mapping to a different concept rather than a mere lack of conceptual possession. We also observe cases where specific collisions only occur part of the time, such as in the case of DALL-E 2 and ZH:ticket (Figure 8 (b)).

Finally, we observed cases of *generic collisions*. For example, DALL-E 2 consistently generates images of desert or Mediterranean scenery when prompted with “eye” in Hebrew (Figure 8 (b)). This pattern shows up across a diverse set of models and prompts. Figure 1 shows how across “dog,” “airplane,” and “face,” DALL-E mega, Stable Diffusion 2, and DALL-E 2 seem to generate vaguely-Israel-looking outdoor scenes regardless of eliciting concept. This is probably reflective of a small sample-size bias in the training data.

### 5.3 Model comparison

Table 2 shows the the use of correctness scores in the *CoCo-CroLa* benchmark to compare the 9 models. As expected, given its multilingual training regimen, AltDiffusion m9 outperforms the other T2I models on average, and in terms of total wins. It is particularly strong relative to the other models in Japanese and Chinese (with the exception of the Chinese-only CogView 2, which is best on Chinese but worst on average overall for both Xc and Wc).

However, despite the strong average performance of AltDiffusion, there’s a lot of room for improvement. For example, its improvements in

terms of JA and HE performance come at a cost of significantly reduced EN and DE performance relative to Stable Diffusion 2, its initialization checkpoint. The *CoCo-CroLa* benchmark can guide future work in adapting T2I models to further multilinguality without losing conceptual coverage on source languages.

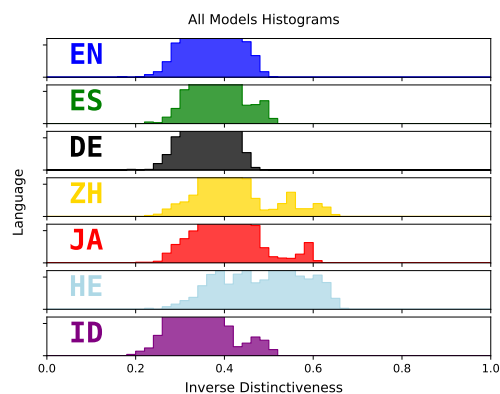


Figure 9: Histograms of the **inverse distinctiveness** scores for all models and all concepts. A model-by-model breakdown is presented in Figure 12.

### 5.4 Distinctiveness captures generic collisions

Figure 9 shows the distribution of the **inverse distinctiveness** score  $D_t$ . On this plot, more rightward probability mass indicates a distribution of concepts for which distinctiveness is **low** relative to a generic sample of images produced by a given model in that language. The four Latin script languages (EN, ES, DE, ID) exhibit the lowest inverse distinctiveness, and are thus the least prone to producing generic failure images. Hebrew is an outlier for having high concept-level  $D_t$ . Across many models. This is probably because it is the most low-resource language in our list by far, and doesn’t benefit from script sharing.

	EN	ES	JA		EN	ES	JA
DALL-E Mega							
	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc
	bird 0.741 27 ✓	bird 0.739 27 ✓	bird 0.704 26 ✓	dog 0.746 26 ✓	dog 0.712 27 ✓	dog 0.298 19 ✗	
	keyboard 0.824 28 ✓	keyboard 0.801 29 ✓	keyboard 0.346 18 ✗	fire 0.938 27 ✓	fire 0.926 27 ✓	fire 0.247 19 ✗	
bird 0.741 27 ✓	bird 0.739 27 ✓	bird 0.704 26 ✓	moon 0.868 29 ✓	moon 0.864 28 ✓	moon 0.269 23 ✗		
AltDiffusion							
	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc
	bird 0.655 26 ✓	bird 0.646 26 ✓	bird 0.655 26 ✓	dog 0.702 26 ✓	dog 0.643 26 ✓	dog 0.677 26 ✓	
	keyboard 0.491 27 ✓	keyboard 0.489 26 ✓	keyboard 0.462 24 ✗	fire 0.669 23 ✓	fire 0.658 23 ✓	fire 0.639 23 ✓	
bird 0.655 26 ✓	bird 0.646 26 ✓	bird 0.655 26 ✓	moon 0.704 27 ✓	moon 0.723 28 ✓	moon 0.607 24 ✓		
Stable Diffusion 2							
	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc	Concept Xc Wc
	bird 0.726 27 ✓	bird 0.697 26 ✓	bird 0.655 26 ✓	dog 0.748 26 ✓	dog 0.712 26 ✓	dog 0.582 25 ✓	
	keyboard 0.837 29 ✓	keyboard 0.789 29 ✓	keyboard 0.797 29 ✓	fire 0.775 25 ✓	fire 0.620 23 ✓	fire 0.292 20 ✗	
bird 0.726 27 ✓	bird 0.697 26 ✓	bird 0.655 26 ✓	moon 0.756 28 ✓	moon 0.763 29 ✓	moon 0.282 19 ✗		

(a) “a bird using a keyboard in the snow”

(b) “a dog made of fire standing on the moon”

Figure 10: Cross-model analysis of more complicated, creative prompts combining concepts including “snow,” “keyboard,” “bird,” “dog,” “fire,” and “moon.” We find that **if a model is found to not possess a concept, it will not be able to produce more complicated prompts including the concept.** This validates *CoCo-CroLa* as an efficient way to capture an overview of a model’s generalization capabilities.

## 5.5 Ranking concepts by Xc

For a given model and language, *CoCo-CroLa* can be used as a concept-level analysis tool. For example, by performing the same ranking over a specific (model, language) pair, we can find the most well-covered and poorly-covered concepts for that pair. For all models and languages, an interactive ranking demo based on ascending and descending Xc and Wc is available at [saxon.me/coco-crola/](https://saxon.me/coco-crola/). For example, we found “snow” to be a concept possessed in EN and ES for DALL-E Mega, AltDiffusion, and Stable Diffusion, but only possessed in JA by Stable Diffusion 2. A similar situation holds for “dog” and “fire,” with AltDiffusion.

## 5.6 Concept possession as a proxy

In this section we will discuss how **an inability for a model to use some concept in complex, creative prompts is implied by our detected non-possession of said concept**, thereby validating the *CoCo-CroLa* atomic evaluation paradigm.

To investigate this we manually translated two creative prompts including concepts found to be differentially present in DALL-E Mega, AltDiffusion, and Stable Diffusion subsection 5.5 from English

into Spanish and Japanese. The prompts were: “a bird using a keyboard in the snow,” (ES: “un pájaro usando un teclado en la nieve,” JA: “雪にキーボードを使っている鳥”) and “a dog made of fire standing on the moon,” (ES: “un perro hecho de fuego pisando en la luna,” JA: “火でできた犬が月に立っている”).

Figure 10 clearly shows that, using thresholds for non-possession of  $Xc < 0.5$  and  $Wc < 25$ , **if a concept is not possessed by a model according to *CoCo-CroLa*, it will be unable to successfully generate creative images containing it.**

However, other capabilities including compositionality and perhaps a sort of *verb-level conceptual possession* are probably also required in order to make the converse (possession implies capability to generate creatively) to be true—yet we have no method to capture such possession. This is a promising direction for future work.

Thus, concept-level coverage in a model can be used as a proxy for generalization capabilities to more complex prompts containing the concept, at least in the case of tangible noun concepts. This is good news, as it **enables assessment of the infinite space of creative prompts from a finite, constrained set of atomic concepts.**



## 6 Conclusion

Multilingual analysis of text-to-image models is desirable both to improve the multicultural accessibility of T2I systems and to deepen our understanding of their semantic capabilities and weaknesses. Analyzing a model’s *conceptual coverage* is a simple and straightforward way to do this.

We demonstrated that these concepts are core building blocks for producing impressive images, and that analyzing the concepts is a useful proxy for assessing the creative capabilities of T2I models.

Our technique, *CoCo-CroLa* is a first step toward further work in this domain. Following our recipe, larger benchmarks containing more languages and concepts can be built easily.

### Limitations

The *CoCo-CroLa* benchmark generating procedure is intended to yield multilingual evaluations that can be scaled to even larger sets of concepts and languages without experienced annotators. In the interests of both concept and language quantity scale, we opted for an automated procedure which leverages machine translation systems, which can introduce translation errors. Furthermore, variation in the nuance or normative meaning of concepts, particularly culturally contested ones such as “face,” (Engelmann et al., 2022) “person,” or “man” will inevitably drive some variance in expected outputs by users across language communities. This cultural variation will place an unavoidable upper bound on the performance of inherently cross-cultural benchmarks such as *CoCo-CroLa*.

Additionally, typological variation between languages can introduce complications in applying our framework. For example, while simple template filling for prompting is straightforward in Chinese, which requires no word-dependent articles, in English phonological properties of the word govern the preceding article, and in Spanish and German grammatical gender do the same. Hebrew has gendered nouns, adjectives, and verbs but not articles, on the other hand. Overall, it appears that these have limited influence as grammaticality isn’t a crucial feature in the prediction of image tokens performed in T2I models, Appendix B.

While doing so aids in the scalability of the approach, using CLIP as a feature extractor for computing the metrics, particularly correctness  $X_c$  and  $W_c$ , potentially introduces biases due to the English-primary data that CLIP is pretrained on.

Future work could test this hypothesis by comparing the performance of *CoCo-CroLa*’s CLIP-based features with  $X_c$  computed using Inception features (as in FID) (Chong and Forsyth, 2020) or with dedicated concept-level purpose-trained classifiers.

### Ethics Statement

Images of human faces are generated by our model. To mitigate the minor risk of resemblance to real people, we have downsampled all images. However, we believe this risk is mitigated by the lack of personal names in the querying data. Furthermore, we believe demonstrating that human faces are generated and under which conditions they are is important for documentation of bias (Paullada et al., 2021) and harm risks in these models.

License information is provided in our project page ([saxon.me/coco-crola](http://saxon.me/coco-crola)) and project repository ([github:michaelsaxon/CoCoCroLa](https://github.com/michaelsaxon/CoCoCroLa)).

### Acknowledgements

This work was supported in part by the National Science Foundation Graduate Research Fellowship under Grant No. 1650114. This work was further supported by the National Science Foundation under Grant No. 2048122. We would also like to thank the Robert N. Noyce Trust for their generous gift to the University of California via the Noyce Initiative. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the sponsors.

### References

- Federico Bianchi, Pratyusha Kalluri, Esin Durmus, Faisal Ladhak, Myra Cheng, Debora Nozza, Tatsunori Hashimoto, Dan Jurafsky, James Zou, and Aylin Caliskan. 2022. *Easily accessible text-to-image generation amplifies demographic stereotypes at large scale*. *ArXiv preprint*, abs/2211.03759.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. *Language models are few-shot learners*. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

- Zhongzhi Chen, Guang Liu, Bo-Wen Zhang, Fulong Ye, Qinghong Yang, and Ledell Wu. 2022. [Altclip: Altering the language encoder in clip for extended language capabilities](#). *ArXiv preprint*, abs/2211.06679.
- Hyundong Cho, Chinnadhurai Sankar, Christopher Lin, Kaushik Ram Sadagopan, Shahin Shayandeh, Asli Celikyilmaz, Jonathan May, and Ahmad Beirami. 2021. [Checkdst: Measuring real-world generalization of dialogue state tracking performance](#). *ArXiv preprint*, abs/2112.08321.
- Jaemin Cho, Abhay Zala, and Mohit Bansal. 2022. [Dall-eval: Probing the reasoning skills and social biases of text-to-image generative transformers](#).
- Min Jin Chong and David Forsyth. 2020. Effectively unbiased fid and inception score and where to find them. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6070–6079.
- Christopher Clark, Mark Yatskar, and Luke Zettlemoyer. 2019. [Don’t take the easy way out: Ensemble based methods for avoiding known dataset biases](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4069–4082, Hong Kong, China. Association for Computational Linguistics.
- Boris Dayma, Suraj Patil, Pedro Cuenca, Khalid Saifullah, Tanishq Abraham, Phuc Le Khac, Luke Melas, and Ritobrata Ghosh. 2021. [Dall-e mini](#).
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, and Jie Tang. 2021. [Cogview: Mastering text-to-image generation via transformers](#). In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 19822–19835.
- Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith. 2020. [Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping](#). *ArXiv preprint*, abs/2002.06305.
- Severin Engelmann, Chiara Ullstein, Orestis Papakyriakopoulos, and Jens Grossklags. 2022. [What people think ai should infer from faces](#). In *2022 ACM Conference on Fairness, Accountability, and Transparency, FAccT ’22*, page 128–141, New York, NY, USA. Association for Computing Machinery.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2021. [The pile: An 800gb dataset of diverse text for language modeling](#). *ArXiv preprint*, abs/2101.00027.
- Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2021. Datasheets for datasets. *Communications of the ACM*, 64(12):86–92.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144.
- Jeremy Howard and Sebastian Ruder. 2018. [Universal language model fine-tuning for text classification](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 328–339, Melbourne, Australia. Association for Computational Linguistics.
- Ben Hutchinson, Jason Baldridge, and Vinodkumar Prabhakaran. 2022. Underspecification in scene description-to-depiction tasks. *AACL*.
- Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V. Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. [Scaling up visual and vision-language representation learning with noisy text supervision](#). In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 4904–4916. PMLR.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. [Dense passage retrieval for open-domain question answering](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online. Association for Computational Linguistics.
- Diederik P. Kingma and Max Welling. 2014. [Auto-encoding variational bayes](#). In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*.
- Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images.
- Mingi Kwon, Jaeseok Jeong, and Youngjung Uh. 2022. [Diffusion models already have a semantic latent space](#). *ArXiv preprint*, abs/2210.10960.

- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. [RACE: Large-scale ReAding comprehension dataset from examinations](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 785–794, Copenhagen, Denmark. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#). *ArXiv preprint*, abs/1907.11692.
- Elman Mansimov, Emilio Parisotto, Lei Jimmy Ba, and Ruslan Salakhutdinov. 2016. [Generating images from captions with attention](#). In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*.
- Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 220–229.
- Roberto Navigli and Simone Paolo Ponzetto. 2010. [BabelNet: Building a very large multilingual semantic network](#). In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 216–225, Uppsala, Sweden. Association for Computational Linguistics.
- Amandalynne Paullada, Inioluwa Deborah Raji, Emily M. Bender, Emily Denton, and Alex Hanna. 2021. [Data and its \(dis\)contents: A survey of dataset development and use in machine learning research](#). *Patterns*, 2(11):100336.
- VU Prabhu and A Birhane. 2020. [Large datasets: A pyrrhic win for computer vision](#). *ArXiv preprint*, abs/2006.16923.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. [Learning transferable visual models from natural language supervision](#). In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR.
- Deborah Raji, Emily Denton, Emily M. Bender, Alex Hanna, and Amandalynne Paullada. 2021. [Ai and the everything in the whole wide world benchmark](#). In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1. Curran.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. [Know what you don’t know: Unanswerable questions for SQuAD](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 784–789, Melbourne, Australia. Association for Computational Linguistics.
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. [Hierarchical text-conditional image generation with clip latents](#). *ArXiv preprint*, abs/2204.06125.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. [Zero-shot text-to-image generation](#). In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 8821–8831. PMLR.
- Scott E. Reed, Zeynep Akata, Xinchun Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. 2016. [Generative adversarial text to image synthesis](#). In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 1060–1069. JMLR.org.
- Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. [Beyond accuracy: Behavioral testing of NLP models with CheckList](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. [High-resolution image synthesis with latent diffusion models](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695.
- Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. 2022. [Photorealistic text-to-image diffusion models with deep language understanding](#). 35:36479–36494.
- Michael Saxon, Xinyi Wang, Wenda Xu, and William Yang Wang. 2021. [Peco: Examining single sentence label leakage in natural language inference datasets through progressive evaluation of cluster outliers](#). *ArXiv preprint*, abs/2112.09237.

Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. [Laion-400m: Open dataset of clip-filtered 400 million image-text pairs](#). *ArXiv preprint*, abs/2111.02114.

Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. 2019. [The woman worked as a babysitter: On biases in language generation](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3407–3412, Hong Kong, China. Association for Computational Linguistics.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. [Beyond the imitation game: Quantifying and extrapolating the capabilities of language models](#). *ArXiv preprint*, abs/2206.04615.

Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. 2022. [Challenging big-bench tasks and whether chain-of-thought can solve them](#). *ArXiv preprint*, abs/2210.09261.

Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2019. [Superglue: A multi-task benchmark and analysis platform for natural language understanding](#). *Advances in Neural Information Processing Systems*, 32:3261–3275.

Wenda Xu, Yi-Lin Tuan, Yujie Lu, Michael Saxon, Lei Li, and William Yang Wang. 2022. [Not all errors are equal: Learning text generation metrics using stratified error synthesis](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 6559–6574, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

## A Details on producing the concept list

**Source language term lists.** We first produce a list of English nouns by collating words in term frequency lists extracted from TV closed captions and contemporary fiction novels from Wiktionary<sup>4</sup>, and filter for the 2000 most frequent words in this combined list, and augment it with class label names from CIFAR100 (Krizhevsky et al., 2009).

**Finding good translations.** We feed the list English words into a custom translation pipeline, which simultaneously queries BabelNet (Navigli and Ponzetto, 2010), and an ensemble of four commercial translation systems: Google Translate, Bing Translate, Baidu Translate, and iTranslate<sup>5</sup>.

In response to an English query, the BabelNet API returns a collection of “SynSets,” subgraphs of a combined multilingual word and entity graphs centered on a node the query word maps to (see Figure 5 for examples). Each subgraph links to multiple other nodes, containing terms in both the source language and the target language. These edges can represent, for example, the titles of Wikipedia articles in different language editions of Wikipedia that are marked as being equivalent, thus ensuring that by checking against SynSet edges, a degree of human validation is included automatically. The synset also contains information about whether a given word is a noun. If it is not a noun, the candidate concept is discarded.

To choose the best translation from those edges, the returned translations into the target languages of the English term from the commercial translation services are *melded* by first sorting all returns by number of languages in the return query (in the case that one translation service covers more languages than others), and filling in missing translations by prioritizing alignment in the shared language translations. If any target language is missing a word for a concept at the conclusion of this process, that concept is discarded from the final list.

**Post-filtering.** Once a list of melded translations from the commercial service is returned, each row in the candidate aligned concept list is checked against the corresponding BabelNet SynSets to ensure each translation is present as a connected node, for pseudo-human evaluation. At the end of this process, a list of approximately 250 concepts is re-

<sup>4</sup>[en.wiktionary.org/wiki/Wiktionary:Frequency\\_lists/Contemporary\\_fiction,.../TV/2006/1-1000](https://en.wiktionary.org/wiki/Wiktionary:Frequency_lists/Contemporary_fiction,.../TV/2006/1-1000)

<sup>5</sup>Using the [translators PyPi package](#).

turned. Finally, we manually remove terms that are verb-noun collisions (e.g. hike) to ensure this ambiguity didn't drive any poor translations. The final list for *CoCo-CroLa v1.0* contains 193 concepts.

### A.1 Full concept list

The (English) concepts are: eye, hand, head, smile, face, room, door, girl, person, man, love, watch, arm, hair, mother, car, mom, dad, table, phone, father, grin, mouth, kid, family, finger, world, shirt, ground, sister, chair, kitchen, woman, beer, hill, metal, hotel, princess, bench, detail, bird, cigarette, history, plastic, pizza, airplane, male, backpack, judge, dragon, sea, bike, female, garden, meal, toy, ship, flame, tail, library, weapon, cd, rope, cafeteria, porch, queen, duck, lake, television, boat, tent, roof, ticket, cop, milk, soldier, tank, thigh, belt, sandwich, bullet, teenager, apple, wine, supply, captain, cheese, feather, mask, prince, beaver, seal, stingray, shark, rose, bottle, mushroom, orange, pear, pepper, keyboard, lamp, telephone, couch, bee, beetle, butterfly, caterpillar, cockroach, tiger, wolf, bridge, castle, house, road, cloud, forest, mountain, camel, chimp, kangaroo, fox, raccoon, lobster, spider, worm, baby, crocodile, lizard, dinosaur, snake, turtle, hamster, rabbit, squirrel, tree, bicycle, train, tractor, jump, men, moon, clothes, neck, fire, tire, teacher, movie, dog, ring, eyebrow, sun, tall, doctor, sky, apartment, shoe, rock, daughter, girlfriend, bar, ball, hallway, tv, teeth, police, field, wife, brain, pants, tongue, cup, computer, bottom, bell, aunt, clock, suit, plate, chocolate, snow, guitar, truck, church, husband, van, blanket, bowl, mama, cookie, hat, monster, ceiling.

All translations are provided at the repository and demo page ([saxon.me/coco-crola](https://saxon.me/coco-crola)).

### B Validating the prompt templates

As mentioned in subsection 4.2, the simple template-based approach to generating prompts for concepts leads to the introduction of grammatical errors, e.g. “a photograph of dog.”

However, it is questionable whether small grammatical or logical errors like missing articles matters for high-resourced, well-covered languages like English. After all, the models are clearly able to generate high quality “photograph of dog” pictures without the word “a” in the sentence (Figure 1). But, does the prompt phrasing matter for lower-performance languages in a model? To investigate the impact of prompt phrasing on conceptual cov-

erage, we tested a variety of English, Spanish and Chinese prompt phrasings on the concepts “dog,” “sea,” “airplane,” and “ship” (selected for their wide distribution across the cross-correlation correctness metric range).

For the English prompts, we experimented with including the articles “a,” “the,” “my,” and “an,” as well as using the words “photograph,” “image,” “photo,” and “picture.” For Spanish, we used variations on the phrase “un foto de” (a photo of), including the same set of articles in English “un/una,” “el/la,” “mi,” “tu,” (your) and “nuestra/o” (our). For Chinese, we tried examples that both included and excluded the possessive particle “的” (de), as well as the words “照片” (zhaopian) and “图片” (tupian) for picture/photograph, and including or excluding the prepended phrase “一张” (yi zhang) to create the meaning “one photograph.” We reran the full 193 concept image generations in those three languages for Stable Diffusion 2 and AltDiffusion.

We found limited impact across all of these dimensions. Full details available in our anonymous demo at [saxon.me/coco-crola](https://saxon.me/coco-crola).

### C Additional Plots

Figure 11 shows the language-wise histograms for correctness scores for all nine models.

Figure 12 shows the language-wise histograms for inverse distinctiveness scores for all nine models. On this set of plots, the tendency for Hebrew to be an outlier in terms of inverse distinctiveness (ie, having lots of *generic collisions* for concept failures) is clearly illustrated. However, other noteworthy outliers are DALL-E Mini and Mega performing worse on Chinese and Japanese (possibly script-driven) and CogView 2 having surprisingly low inverse distinctiveness for non-Chinese (non-training) languages in spite of low correctness.

Figure 13 shows the language-wise histograms for the self-consistency scores for all nine models.

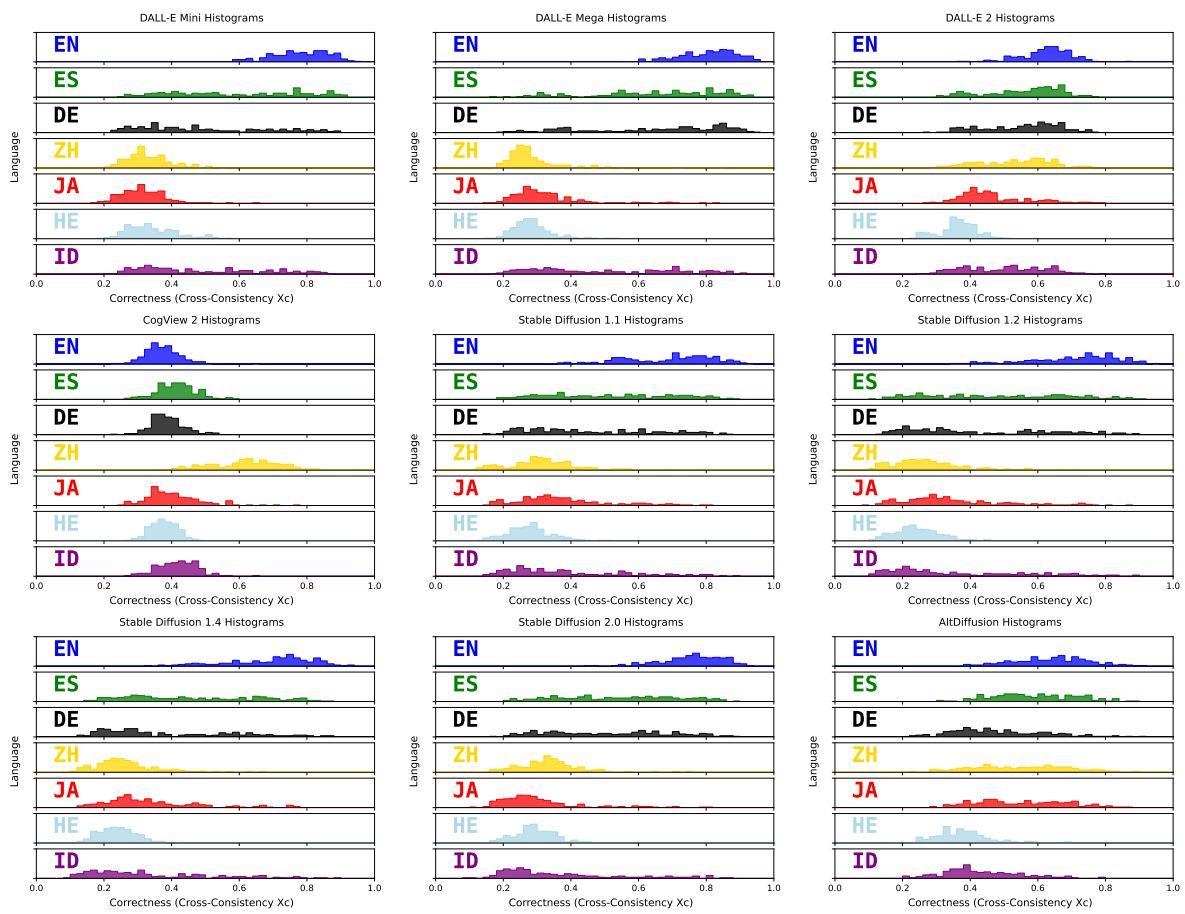


Figure 11: Histograms of the distribution of **correctness** cross-consistency (Xc) for each test language for all assessed models. Rightward probability mass reflects better conceptual coverage.

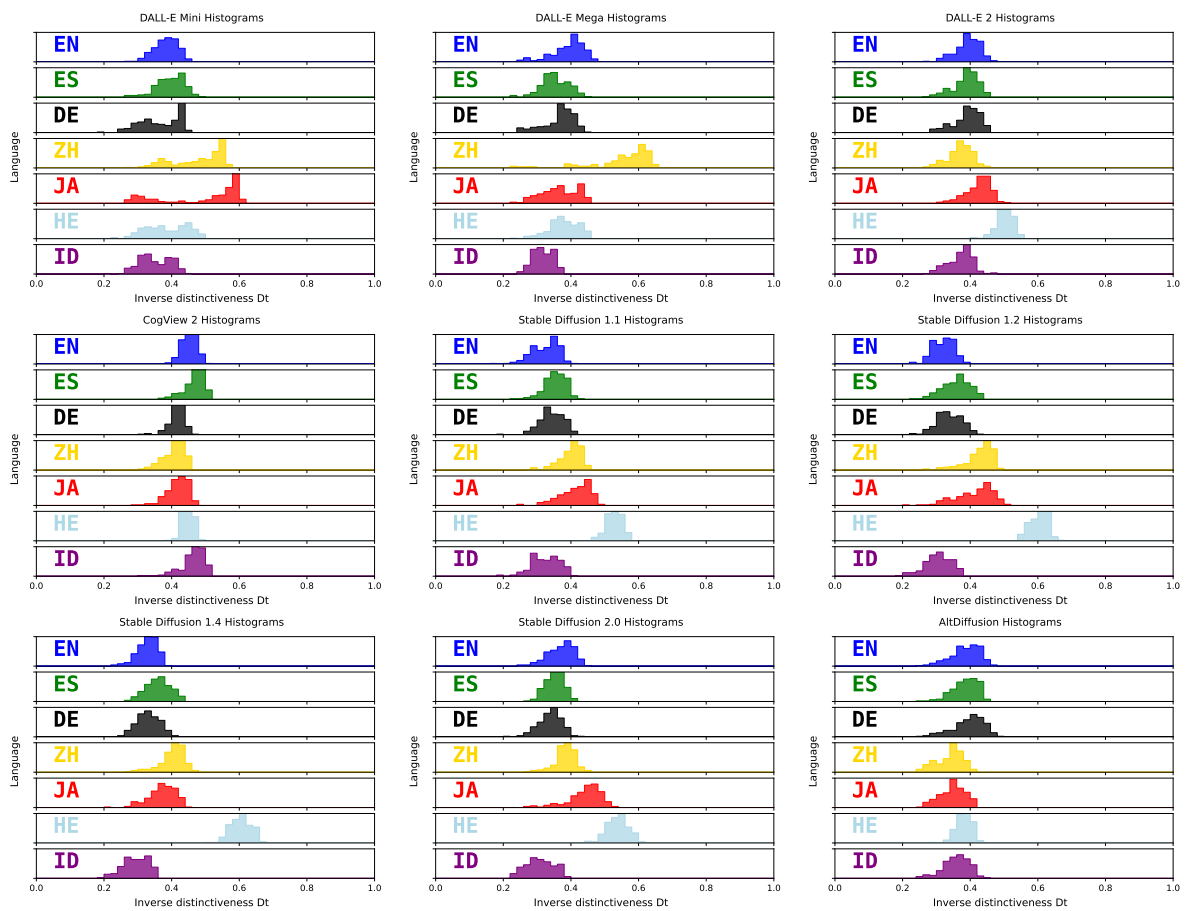


Figure 12: Distribution of **inverse distinctiveness** scores (Dt) for each test language for all models.

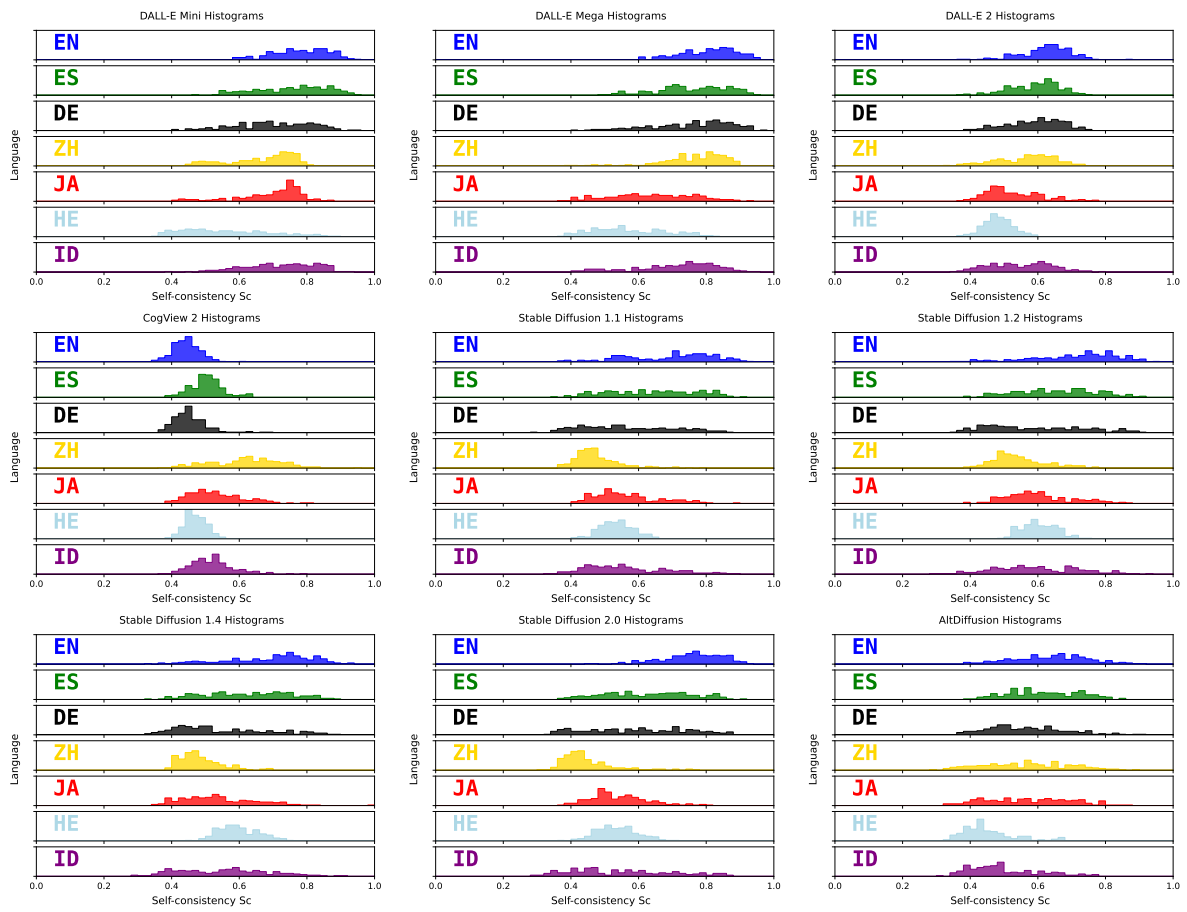


Figure 13: Distribution of **self-consistency** scores (Sc) for each test language for all assessed models.



## ACL 2023 Responsible NLP Checklist

---

### A For every submission:

- A1. Did you describe the limitations of your work?  
*Limitations*
- A2. Did you discuss any potential risks of your work?  
*Ethics Section*
- A3. Do the abstract and introduction summarize the paper's main claims?  
*abs, 1*
- A4. Have you used AI writing assistants when working on this paper?  
*Left blank.*

### B Did you use or create scientific artifacts?

*Throughout, 3,4*

- B1. Did you cite the creators of artifacts you used?  
*Throughout, 4, Appendix A*
- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?  
*Ethics*
- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?  
*4, Appendix A*
- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?  
*No PII risk in common nouns*
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?  
*Throughout*
- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.  
*Left blank.*

### C Did you run computational experiments?

*4*

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?  
*Only used preexisting models, with reference to them for param info*

*The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.*

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

*No training involved*

- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?

5

- C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?

*Appendix A*

**D  Did you use human annotators (e.g., crowdworkers) or research with human participants?**

*Left blank.*

- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?

*No response.*

- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?

*No response.*

- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?

*No response.*

- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

*No response.*

- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

*No response.*