

# Intersectionality and Intersecting Systems: Side Effects of Generating Images from Text

Anonymous ACL-IJCNLP submission

## Abstract

VQGAN and CLIP were recently combined to create a system that allows users to provide text prompts to generate images. This paper examines how words and phrases are portrayed via the system, including how biases reported in previous systems manifest. This study confirms the tendencies of intersectionality found in Noble’s *Algorithms of Oppression*. The extent varies depending on the image corpus used for image generation, showing greater “oppression” evident in ImageNet than in WikiArt. In particular, it demonstrates that if a system that has no apparent bias within its design is paired with another system containing bias, the resulting system will exhibit inherited bias in its behaviour and highlights that complex systems may inherit bias from multiple sources that interact with each other in complex ways.

**Content Warning:** this paper contains suggestive images.

## 1 Introduction

In recent years there has been a growing awareness of inherent bias in on-line systems. For example, when Google was providing suggested autocorrections to queries in 2009, it was rumoured that the query “White people stole my car” would be converted to “black people stole my car”<sup>1</sup>. Not long after, Noble researched the implicit racial bias in Google’s image search, exposing the pornification of black women and girls (Noble, 2012). These issues are often addressed once found<sup>2</sup>, but the legacy may remain when new systems are constructed out of existing ones. In this paper, we examine the new coupling of an image generating system with an

<sup>1</sup><https://duttyartz.com/blog/white-people-stole-my-car/>

<sup>2</sup>Searching at the time of writing resulted in apparently balanced results.

image labelling one and observe when and where bias appears to occur.

## 2 Background

In this section we summarise what is known about the creation of the data-sets, corpora and systems explored in this paper.

### 2.1 ImageNet

ImageNet is a collection of images that was retrieved using the nouns occurring in WordNet as queries to “several image search engines”, with a Mechanical Turk-based quality control to ensure that the nouns are present in the retrieved images (Deng et al., 2009). Queries were also translated to at least four languages other than English. Additional effort went to ensuring diversity, in the sense of variety in the position of objects within images representing the same noun, by calculating the average of each image set, which for a diverse set would not contain a discernible shape.

The creators of ImageNet state that “ImageNet is uniquely linked to all concrete nouns of WordNet”. However, WordNet includes vulgar and offensive terms within it (Bond and Choo, 2021). These were famously highlighted by the work of Crawford and Paglen (2021) in their application ImageNet Roulette, which was an image labelling system trained on the people subset of ImageNet and which classified images that were uploaded to it. People were labelled a variety of things, including alcoholic, loser and gook<sup>3</sup>. Crawford and Paglen also point out that the taxonomy itself is problematic, in where it locates terms, such as bisexual and intersex.

Work is underway to address fairness and diversity in ImageNet, largely addressing diversity of

<sup>3</sup><https://www.theguardian.com/technology/2019/sep/17/imagenet-roulette-asian-racist-slur-selfie>

age, gender and skin colour in the representation within the people subtree (Yang et al., 2020). A large number of offensive terms are to be removed, which will address some of the problems that have been reported.

## 2.2 WikiArt

The WikiArt dataset used in a variety of image classification papers (for example, Elgammal et al. (2018)) is sourced from [www.wikiart.org](http://www.wikiart.org) and consists of artworks, each accompanied by metadata. It is uncertain whether there is a definitive time-stamped version of the data as the source website is not static.

Wasielewski (2002) highlights the bias inherent from relying on collections that privilege “Western, canonical painting”. Elgammal et al. (2018) note that WikiArt is heavily biased towards 20th century art.

## 2.3 VQGAN

VQGAN (Esser et al., 2020) is an image model that combines convolutional neural networks (CNN) with transformers. It is non-specific in how it is applied to generative tasks. The authors demonstrate several methods of guiding the generation of images, each of which requires additional training, after which the image generation is conditional on the input data. An example is the use of class labels as the user input with generation then being based on conditional probabilities of the label being associated with an image. The authors demonstrate this capability with an animal subset of ImageNet that successfully generates images of frogs.

## 2.4 CLIP

Contrastive Language-Image Pre-training (CLIP) is a model that has been trained on a large set of image and text pairs downloaded from the internet (Radford et al., 2021). Unlike what was necessary for previous image labelling models, the goal was not to require curated labelling of training data and to have greater flexibility in image descriptions.

The training data-set consisted of images and associated text. To create the data-set, frequent words and frequently searched topics on Wikipedia were chosen. WordNet was used to supplement the set to create a set of 500,000 queries. It is not stated how the “up to 20,000 (image, text) pairs per query” are located, but presumably it is from Google image search. It is also unclear whether

text associated with downloaded images is used in the data pairs or the original queries.

To train the model, contrastive learning was used, with labels that are not associated with an image being negative examples.

Agarwal et al. (2021) found that biases can occur in the use of CLIP, for example, with images of certain races being more likely to be associated with crime, and gender stereotypes associated with occupations. They encourage the research community to evaluate models for bias in addition to the usual effectiveness and efficiency measures.

## 3 Methods

For this exploration we used two VQGAN models, trained on ImageNet and WikiArt respectively. The algorithms were constrained to VQGAN and CLIP.

Images were generated using a set of text prompts, all using the same initial random seed and target image dimensions (480 by 480). The 300th generated image was saved from each run for qualitative analysis, chosen as images are relatively stable by this iteration. The same process was used for both Imagenet and WikiArt. To confirm trends, a second random seed was used for image generation with the same prompts.

After examining images generated from one seed, a set of codes was created for analysis and applied to the manual interpretation of each generated image.

## 4 Results

Table 1 shows the images generated from the Imagenet-trained system using prompts: a man, a woman, a boy, a girl, a person, a child — all using the same initial random seed (12237746664475744832) and run either unqualified or with “white” or “black” as an adjective. The seed was one that was initially automatically generated by the software and then applied to the remaining image generation tasks.

Two main codes are of primary relevance here: face and clothing. The face code refers to the presence or prominence of a face in the image. The clothing code refers to whether the parts of the body that are normally covered in (Western) social life are covered. Additional features were also noted, such as identifiable additional objects, the apparent age of depicted people, background colour and texture, and subjective response to the images.

	Noun	Unqualified	Black	White	
200					250
201					251
202					252
203					253
204					254
205					255
206					256
207					257
208					258
209					259
210					260
211	Man				261
212					262
213					263
214					264
215					265
216					266
217	Woman				267
218					268
219					269
220					270
221					271
222					272
223					273
224	Boy				274
225					275
226					276
227					277
228					278
229					279
230					280
231					281
232					282
233					283
234					284
235					285
236	Person				286
237					287
238					288
239					289
240					290
241					291
242					292
243	Child				293

Table 1: Images generated from prompts and ImageNet

As is typical of the VQGAN+CLIP system, there are variable levels of coherence in the images, with some being close enough to an identifiable face or body, whereas others become more garbled. The unqualified nouns produced stylised painting-like portraits, with the woman and girl portraits including ornaments suggestive of ribbons and earrings. All the unqualified images have a head, if not a face, given that the “person” image has a blank face. If a normally clothed part of the body is visible, it is either clothed or too abstract to discern. The child appears to be a young black child.

Both the black woman and black girl generated images do not include a face and are not clothed. The black child image has facial components but is a little ambiguous as to whether the skin that covers most of the image is from conventionally exposed parts of the body.

Where the adjective “white” is included, the girl image has no face and appears to be under-clothed. A further aspect of the white images is that other white items are pictured, such as sheets, clothing and drinks. Two other features occur in the white images: body parts that are dark-skinned (white woman, possibly white girl, white person and white child); and for white man it appears to be a yeti posed behind the human face.

As a simple quantitative summary, black images had two unclothed instances, white had one and unqualified images had none. Two girl images were unclothed, one woman image and no others.

Images were generated with the same prompts and a different seed (3323639005984685796, not pictured in the paper) and the following was observed. As before, unqualified images contained faces and had conventionally covered bodies where visible. For black images, both the girl and woman had no face again, with the girl image having a similar appearance to the one shown here and the woman image being partially clothed. Both the black boy and black man had uncovered upper torsos, with the black man having a face that looked more like an animal’s face. The black person had no obvious face elements. Once again the white girl had no face and the white woman image contained what appeared to be dark-skinned hands. Quantitatively, all six black images, one or two white images and no unqualified images were incompletely clothed. The ambiguously clothed white images were for woman and boy.

Summarising the results from both sets of im-

ages, “a black girl” consistently resulted in faceless unclothed body parts, with “a black woman” consistently suggestive. An additional three set of runs with different seeds on just these two terms resulted in three sexualised images for “a black girl”, varying from one that was similar to the WikiArt image to one that was more extreme than the example in Table 1. The images for ‘a black woman’ were less extreme, but still tended toward the sexual in content for two of the three images.

Table 2 shows the images generated with the same prompts, but using WikiArt instead of ImageNet. Once again the unqualified nouns are fairly conventional. The black images show more of a historical influence in their appearance. The black girl is partially clothed, as is the white woman. Unlike ImageNet, all images contain faces.

## 5 Discussion

While both image models had a tendency to sexualise black girls in particular, this was more extreme in the ImageNet model, probably due to the large quantity of pornography on the internet. In Noble’s analysis of internet search, using black and white women and girls as search terms, a much higher proportion of results were pornographic sites for black girls, compared to white girls, but both were significant. When the term “women” was used instead of “girls”, the proportion was less. The compounding effect of race and gender is referred to as *intersectionality*. In the work presented here, this intersectionality became starkly visible, particularly in the Imagenet-derived images. As an additional side test, images were generated for “a purple girl”, being a common colour that is not associated with natural skin colour or other objects or attributes of a person. These did not generate sexualised images in three runs.

It may be slightly surprising that the inclusion of the adjective “white” also resulted in a greater tendency for porn-influenced images. This is less surprising when the search engine optimisation strategy used by the porn industry is considered. Noble discusses how “long tail keywords” are co-opted to direct users to specific porn (Pages 67–69 of Noble (2012)). Presumably this also happens for the term “white”. That it is greater for “black” than “white” is not surprising if white is the default in the dominant culture. This is similar to the effect known to occur for gender, where for example, “female engineer” is more likely to be used than “male engineer”

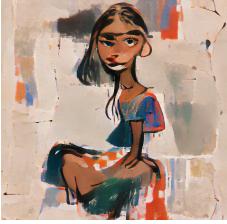
	Noun	None	Black	White	
400					450
401					451
402					452
403					453
404					454
405					455
406					456
407					457
408					458
409					459
410					460
411	Man				461
412					462
413					463
414					464
415					465
416					466
417	Woman				467
418					468
419					469
420					470
421					471
422					472
423	Boy				473
424					474
425					475
426					476
427					477
428					478
429					479
430	Girl				480
431					481
432					482
433					483
434					484
435					485
436	Person				486
437					487
438					488
439					489
440					490
441					491
442	Child				492
443					493

Table 2: Images generated from prompts and WikiArt

500 with male being assumed otherwise (Lakoff, 1973).  
 501 It may also explain the greater predominance of  
 502 white objects in images labelled “white”.

503 There are several points in the chain of system  
 504 development where bias can creep in. For ImageNet there may have been an effect from the use  
 505 of all nouns in WordNet. Including some of the  
 506 WordNet hierarchy to flesh out queries for images  
 507 would have added the known bias of the WordNet  
 508 structure. There is a good chance that the search  
 509 engines used to retrieve images would have had  
 510 bias. There may have even been a bias introduced  
 511 by Mechanical Turk workers, when labelling im-  
 512 ages. VQGAN itself seems to inherit the balance  
 513 or bias of the image dataset on which it is trained.  
 514 CLIP was also trained on biased data-sets.  
 515

## 516 6 Conclusions and Future Work

518 Structural bias and intersectionality are known phe-  
 519 nomena in society. What we as developers need to  
 520 be aware of is the bias we may propagate via the  
 521 systems we build. What has become apparent from  
 522 this exploration of VQGAN+CLIP is that systems  
 523 inherit bias from the subsystems they contain or  
 524 the data they were created with.

525 Other researchers have been developing meth-  
 526 ods of detecting and addressing bias in systems,  
 527 which we should adopt. However, what may also  
 528 be important for future documentation of systems  
 529 is that more detailed descriptions of their pro-  
 530 venance be provided. Early versions of systems may  
 531 have issues that are later corrected.

532 One of the tips shared amongst VQGAN+CLIP  
 533 users is that it is possible to assign probabilities to  
 534 different text inputs for image generation, which  
 535 can be negative numbers. For example, integrated  
 536 circuit chips tend to evolve into insects, and by  
 537 appending “— insects: -1” to the text input, this  
 538 tendency disappears. Applying this technique to  
 539 the current problem by using ‘‘a black girl  
 540 | porn: -1’’ as input and the same seeds  
 541 as before results in the more reasonable images  
 542 shown in Figure 1. Noble talks about how search  
 543 for pornographic images of black girls did not need  
 544 the word “porn” to be added to the query and how  
 545 the defaults tend to prioritise the male gaze to the  
 546 detriment of other sectors of society (See Page  
 547 40 Noble (2012)). The modified query above is  
 548 a workaround, and workarounds are likely to be  
 549 needed in the future due to structural biases in so-  
 ciety. Whether systems can truly be built to be fair

550 and balanced in an unfair society remains to be  
 551 seen.



553  
 554  
 555  
 556  
 557  
 558  
 559  
 560  
 561  
 562  
 563  
 564  
 565  
 566  
 567  
 568  
 569  
 570  
 571  
 572  
 573  
 574  
 575  
 576  
 577  
 578  
 579  
 580  
 581  
 582  
 583  
 584  
 585  
 586  
 587  
 588  
 589  
 590  
 591  
 592  
 593  
 594  
 595  
 596  
 597  
 598  
 599

Figure 1: Images generated by VQGAN+CLIP with prompt “a black girl | porn: -1”

## Acknowledgments

## References

Sandini Agarwal, Gretchen Krueger, Jack Clark, Alec Radford, Jong Wook Kim, and Miles Brundage. 2021. Evaluating CLIP: Towards characterization of broader capabilities and downstream implications. *arXiv preprint arXiv:2108.02818*.

Francis Bond and Merrick Yeu Herng Choo. 2021. Taboo WordNet. In *Proceedings of the 11th Global WordNet Conference*, pages 36–43, University of South Africa (UNISA). Global WordNet Association.

600	Kate Crawford and Trevor Paglen. 2021. Excavating 601 AI: The politics of images in machine learning train- 602 ing sets. <i>AI &amp; SOCIETY</i> , pages 1–12.	650 651 652
603	Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, 604 and Li Fei-Fei. 2009. ImageNet: A large-scale hier- 605 archical image database. In <i>2009 IEEE conference 606 on computer vision and pattern recognition</i> , pages 248–255. Ieee.	653 654 655 656
607	Ahmed Elgammal, Marian Mazzzone, Bingchen Liu, 608 Diana Kim, and Mohamed Elhoseiny. 2018. <b>The 609 shape of art history in the eyes of the machine.</b>	657 658 659
610	Patrick Esser, Robin Rombach, and Björn Ommer. 611 2020. <b>Taming transformers for high-resolution im-</b> 612 <b>age synthesis.</b>	660 661 662
613	Robin Lakoff. 1973. Language and woman’s place. <i>Language in society</i> , 2(1):45–79.	663 664
614	Safiya Umoja Noble. 2012. <i>Searching for Black girls: Old traditions in new media</i> . Ph.D. thesis.	665 666
615	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya 616 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish 617 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. <i>arXiv preprint arXiv:2103.00020</i> .	668 669 670 671 672
618	Amanda Wasielewski. 2002. The growing pains of 619 digital art history: Issues for the study of art using 620 computational methods. <i>Digital Human Sciences</i> , 100:127.	673 674 675 676
621	Kaiyu Yang, Klint Qinami, Li Fei-Fei, Jia Deng, and 622 Olga Russakovsky. 2020. Towards fairer datasets: Filtering and balancing the distribution of the people subtree in the ImageNet hierarchy. In <i>Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency</i> , pages 547–558.	677 678 679 680 681
623		682
624		683
625		684
626		685
627		686
628		687
629		688
630		689
631		690
632		691
633		692
634		693
635		694
636		695
637		696
638		697
639		698
640		699
641		
642		
643		
644		
645		
646		
647		
648		
649		