**Title**

Developing a Culturally Relevant Vocabulary for Talking About Bias in NLP Systems

**Introduction and previous work(s)**

The use of ultra-large language models in natural language processing have recently led to big breakthroughs in previously intractable tasks like Question-Answering and machine translation.[1] However, these models are typically trained on datasets curated from internet resources, including the Reddit Corpus, Common Crawl, and Wikipedia, which have been shown to exhibit racial and gender bias.[2] Many attempts to de-bias models have been proposed, and bias tasks such as WEAT and Winobias exist to try and quantify the amount of bias these models exhibit.[3,4] Although we have many methods that quantify bias in natural language models, the qualitative labels used to measure this bias have not been as rigorously studied. For example, in Bolukbasi et al, gendered words are calculated from word embeddings, while stimuli in WEAT include words like "diamond" and "paradise" for pleasantly associated words, and words like "crash" and "accident" for unpleasantly associated words. These supposedly biased words a casual internet user likely wouldn't associate with any particular group identity. Meanwhile, stimuli words seen in the most egregious examples of bias are often not present in most bias datasets. (Figure 1) By not using culturally-relevant terms for bias measurement, the bias measured in these quantitative tests will by definition themselves be biased.

**Proposed Project**

I propose to then build a list of culturally relevant bias terms, and will propose a new method of testing bias in generative language models. First, I will evaluate the current cultural relevance of the proposed bias stimuli from the Winobias and WEAT datasets. To do so, I will conduct a user study on Mechanical Turk with 500 participants to see if and how strongly people associate commonly-used stimuli from these bias datasets with different group identities. I will then scrape examples of explicit bias in language models from Twitter and the Internet, similar to those found in Figure 1. I will then compare the word embedding distances of these scraped words using the methods of Bolukbasi et al to the distances of the standard stimuli words from WEAT and Winobias, and I will conduct another user study with the scraped words. Finally, I will re-perform the analyses of WEAT and Winobias using only *relevant* bias stimuli, which I define as words that have both a biased word vector representation and are labeled as culturally relevant in the aforementioned user studies. I hope to submit my work to ACL 2021.

**Specific Milestones**

- Run user study of current bias stimuli - 2 weeks to complete
- Scrape culturally-relevant bias examples from the Internet and Twitter - 1 month
- Compare to current word embeddings and run new user study - 2 weeks
- Re-analyze WEAT, Winograd, Winobias with relevant stimuli - 3 weeks

**References**

[1] Devlin et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. Arxiv, 2019.

[2] Bolukbasi et al. Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings. *NeurIPS*, 2016.

[3] Caliskan et al. Semantics derived automatically from language corpora contain human-like biases. *Science*, 2017

[4] Zhao et al. Gender Bias in Coreference Resolution: Evaluation and Debiasing Methods. *NAACL,* 2018.

Figure 1: example of racial bias from GPT-3

**Budget**

| Item | Cost |
| --- | --- |
| 2 MTurk User studies | $500 |
| Cloud compute | $500 |
| ACL conference fees | $500 |

**Project Title:** Creating a dataset to study the evolution of Spatial Apartheid in South Africa

**Introduction and previous work(s) (200 words max):**

Aerial images of residential areas in South Africa show the clear legacy of spatial apartheid, a former policy of political and economic discrimination against non-European groups, with completely segregated neighbourhoods of townships next to gated wealthy neighbourhoods. This proposal seeks to work on the first publicly available computer vision dataset to study the evolution of spatial apartheid in South Africa. The dataset will consist of satellite images of South Africa from 2006-2017, some of which will be annotated by neighborhood classes at the pixel level. While there are several freely available datasets for the broader task of land cover classification, DeepGlobe [1], UC Merced [2], etc., most of them have been created for the developed world [3, 4]. While datasets such as [1], were created using images from developing countries (Afghanistan, Thailand, Indonesia, India), neighbourhoods in these countries also have very different characteristics to those in South Africa [5]. On the other hand, publicly available datasets denoting land-use in the African continent usually only distinguish between informal settlements and everything else, and only cover a single city due to cost constraints [6].

**Proposed project (200 words max):**
We plan on creating a dataset consisting of annotated satellite images of South Africa from 2006-2017, labeled by neighborhood type. To create our dataset, we plan on using satellite images covering the entire country of South Africa from the South African National Space Agency [7]. To associate pixels in the images with the types of neighbourhoods they depict, we turn to the Enumeration Areas (EAs) dataset created in 2011 by Statistics South Africa (Stats SA) [8]–a government agency responsible for conducting the census, and the Geo-referenced buildings dataset locating all buildings in South Africa. We may also need to augment our annotations of townships with expert annotators. While building our dataset, we plan to use an iterative approach, first training a U-Net [9] segmentation model, examining the results to see if additional annotations are needed and iterating accordingly. We plan to evaluate the characteristics of our dataset by training and evaluating state-of-the-art segmentation models [10] and performing experiments examining our ability to distinguish between wealthy and non wealthy neighborhoods.

**Specific milestones (3+ milestones recommended, ~100 words max):**
Since creating datasets like this may take months or even years, we do not expect to finish the entire dataset creation process in 3 months, but have carved out specific milestones below which would help us gather an initial database for one province.
- Month 1: Gather the annotated datasets outlined above and align polygons with locations for 1 province (Gauteng).
- Month 2: Recruit expert annotators to gather additional annotations for townships in Gauteng.
- Month 2: Start training U-Net model on Gauteng, start analyzing results.
- Month 3: Start iterating on gathering additional annotations for Gauteng based on results of U-Net model

**Budget (Does not count towards page limit)**

| Cloud compute costs | $750 |
|---|---|
| Pay expert annotators | $500 |
| Unforeseen costs | $250 |

## References (Does not count towards page limit)

1. Ilke Demir, Krzysztof Koperski, David Lindenbaum, Guan Pang, Jing Huang, Saikat Basu, Forest Hughes, Devis Tuia, and Ramesh Raska. Deepglobe 2018: A challenge to parse the earth through satellite images. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Work- shops (CVPRW)*, pages 172–17209. IEEE, 2018.
2. Yi Yang and Shawn Newsam. Bag-of-visual-words and spa- tial extensions for land-use classification. In *Proceedings of the 18th SIGSPATIAL international conference on advances in geographic information systems*, pages 270–279. ACM, 2010.
3. Caleb Robinson, Le Hou, Kolya Malkin, Rachel Soobit- sky, Jacob Czawlytko, Bistra Dilkina, and Nebojsa Jojic. Large scale high-resolution land cover mapping with multi- resolution data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12726– 12735, 2019.
4. Collin Homer, Jon Dewitz, Limin Yang, Suming Jin, Patrick Danielson, George Xian, John Coulston, Nathaniel Herold, James Wickham, and Kevin Megown. Completion of the 2011 national land cover database for the conterminous united states–representing a decade of land cover change in- formation. *Photogrammetric Engineering & Remote Sens- ing*, 81(5):345–354, 2015.
5. Rizwan Ahmed Ansari, Rakesh Malhotra, and Krishna Mo- han Buddhiraju. Identifying informal settlements using con- tourlet assisted deep learning. *Sensors*, 20(9):2733, 2020.
6. Nicholus Mboga, Claudio Persello, John Ray Bergado, and Alfred Stein. Detection of informal settlements from vhr images using convolutional neural networks. *Remote sensing*, 9(11):1106, 2017.
7. SANSA. South African National Space Agency. https: //www.sansa.org.za/, 2019.
8. Statistics SA. Statistics South Africa. http://www. statssa.gov.za/?page_id=3917, 2019.
9. Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U- net: Convolutional networks for biomedical image segmen- tation. In *International Conference on Medical image com- puting and computer-assisted intervention*, pages 234–241. Springer, 2015.
10. Huang, Gao, et al. "Densely connected convolutional networks." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017.

Building information models(BIMs) from 3D point cloud scans

**Introduction and previous work(s)**
In the Architectural/Engineering and Construction (AEC) field, Building Information Models(BIMs) are used to represent the structure and elements of a building with associated topological relationships. BIMs are key for the design and maintenance of buildings and structures, but are difficult to obtain for existing buildings. Previous works have leveraged neural networks to provide building information such as floorplans from 3D point cloud scans[1]. These works produce simplified building information and are limited by lack of available data. A recent effort to automate Scan-to-BIM processes provides an extensive benchmark with 20 buildings, 49 point clouds, and 2D building models for each floor[2]. This benchmark makes it possible to develop and extensively evaluate automated Scan-to-BIM. I propose to implement an automated method to generate BIM data from 3D point cloud scans with a focus on floorplan estimation. Given a 3D scan of a building in the form of a 3D point cloud, the goal is to generate the floorplan of the building including rooms, walls, doors, and windows.

**Method**
For this project, I will leverage the dataset and metrics from the '1st Workshop and Challenge on Computer Vision in the Built Environment for the Design, Construction and Operation of Buildings'[2]. I will first implement a baseline method that extracts the floorplan of a building based on handcrafted features such as histogram density similar to [3]. I will then build on the existing method from Liu et al.[1] to design and train a neural network model that predicts floorplan information given a 3D point cloud scan of a building or room. This model will leverage PointNet[4] to process the input cloud, and optionally ResNet[5] to process images from the building as additional information for floorplan estimation. I will evaluate and compare both methods using metrics including - IoU of the floor area room corner accuracy, orientation, and length of walls, doors, and windows, and room connectivity as defined in [2]. I will explore potential network architecture variations to improve the performance of my model. This work will be submitted as an entry for the '1st Workshop and Challenge on Computer Vision in the Built Environment for the Design, Construction and Operation of Buildings' at CVPR 2021[2].

**Milestones**
  ● Download and pre-process 3D point cloud scans from benchmark. Build data loader and precompute building features - 2 weeks
  ● Design and evaluate baseline model based on handcrafted features - 3 weeks
  ● Implement and evaluate neural network model - 4 weeks
  ● Explore network architecture variations and submit best model as entry to CVPR challenge - 2weeks

**References**
[1]Liu, Chen, Jiaye Wu, and Yasutaka Furukawa. "Floornet: A unified framework for floorplan reconstruction from 3d scans." *Proceedings of the European Conference on Computer Vision (ECCV)*. 2018.
[2]*Computer Vision in the Built Environment*, cv4aec.github.io/.
[3] Armeni, Iro, et al. "3d semantic parsing of large-scale indoor spaces." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016.
[4]Qi, Charles R., et al. "Pointnet: Deep learning on point sets for 3d classification and segmentation." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017.
[5]He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.

**Budget and justification**

I plan to use AWS Cloud Services for storage and computational resources needed for storing the project's data and training my neural network models. I also plan to attend to attend and submit my work to the 1st Workshop and Challenge on Computer Vision in the Built Environment for the Design, Construction and Operation of Buildings' at CVPR 2021

CVPR 2021 Non-member Student passport: **$550**
AWS Cloud compute: **$950**