Responsible NLP Checklist

Paper title: Seeing Through Words, Speaking Through Pixels: Deep Representational Alignment Between Vision and Language Models

Authors: Zoe Wanying He, Sean Trott, Meenakshi Khosla

How to read the checklist symbols:	
the authors responded 'yes'	
the authors responded 'no'	
the authors indicated that the question does not apply to their work	
the authors did not respond to the checkbox question	
For background on the checklist and guidance provided to the authors, see the Responsible NLP Checklist page at ACL Rolling Review.	:

✓ A. Questions mandatory for all submissions.

- A1. Did you describe the limitations of your work? *This paper has a Limitations section.*
- We did not include a discussion on potential risks because our work is a purely offline and descriptive analysis of pre-trained vision and language models internal representations. We do not introduce any new model, dataset, or deployment scenario that directly interacts with users or makes decisions affecting people. Therefore, we believe that there is no immediate ethical harms unique to our methodology.
- **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models)
 - ☑ B1. Did you cite the creators of artifacts you used?

 We cited the creators of the artifacts we use in section 2, Methods.
 - ☑ B2. Did you discuss the license or terms for use and/or distribution of any artifacts?

 Yes. We will release our code in a public GitHub repository under the MIT license (see Appendix: Code Availability). This provides clear terms for use, modification, and redistribution.
 - 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)?
 - Yes. All datasets and models used (MS-COCO, Pick-A-Pic, Flickr8k, BLOOM, OpenLLaMA, Qwen, Phi-3, SmolLM, DINOv2, ViT) were employed under their intended research-use conditions. The artifacts we release (code) are intended for research purposes only and are consistent with the original access conditions of the resources used.
 - 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?
 - In our paper, we only used the public MS-COCO, Pick-A-Pic, and Flickr8k images, and also the synthetic images generated by diffusion models. Therefore, the data should not contain any personally identifying information or offensive content.

- ☑ B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.? Section 2, Methods; Appendix ☑ 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? Section 2, Methods; Appendix **☑** C. Did you run computational experiments? 🛮 C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used? We ran alignment analyses with inference-only on off-the-shelf model, and we used the Gemini API for caption generation. W didnt track GPU time or billing specifics in the manuscript. 2 C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values? Section 2, Methods; Appendix 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? Section 3: Results; Appendix C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation, such as NLTK, SpaCy, ROUGE, etc.), did you report the implementation, model, and parameter settings Section 2, Methods; Appendix **D.** Did you use human annotators (e.g., crowdworkers) or research with human subjects? 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.? (left blank) (e.g., country of residence)? (left blank)
 - 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
 - D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating (e.g., did your instructions explain how the data would be used)? (left blank)
 - □ D4. Was the data collection protocol approved (or determined exempt) by an ethics review board? (left blank)
 - D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data? (left blank)

E. Did you use AI assistants (e.g., ChatGPT, Copilot) in your research, coding, or writing?

■ E1. If you used AI assistants, did you include information about their use? AI only assisted the writing with language including revising the wordings and syntax based on my own words from the draft. AI also helped with Latex formatting. Based on the policy, such form of AI assistance does not need to be disclosed.