

Is Multimodal Vision Supervision Beneficial to Language?

Avinash Madasu

Department of Computer Science
UNC Chapel Hill, USA

avinashmadasu17@gmail.com

Vasudev Lal

Cognitive Computing Research
Intel Labs, USA

vasudev.lal@intel.com

Abstract

Vision (image & video) - Language (VL) pre-training is the recent popular paradigm that achieved state-of-the-art results on multi-modal tasks like image-retrieval, video-retrieval, visual question answering etc. These models are trained in an unsupervised way and greatly benefit from the complementary modality supervision. In this paper, we explore if the language representations trained using vision supervision perform better than vanilla language representations on Natural Language Understanding and common-sense reasoning benchmarks. We experiment with a diverse set of image-text models such as ALBEF, BLIP, METER and video-text models like ALPRO, Frozen in Time, VIOLET. We compare the performance of language representations of stand-alone text encoders of these models to the language representations of text encoders learnt through vision supervision. Our experiments suggest that vanilla language representations show superior performance on most of the tasks. These results shed light on the current drawbacks of the vision-language models. The code is available at <https://github.com/avinashsai/MML>

1. Introduction

Vision-language (VL) pre-training [1, 5, 11, 12, 22] has shown tremendous success in the areas of image-text retrieval [11, 12], visual question answering [4, 31], video retrieval [1, 5, 17, 18]. These models benefit from the mutual supervision of vision and language leading to the superior results on multi-modal tasks. So, the natural question arises: “Are vision supervised language representations beneficial compared to vanilla language representations on Natural Language Understanding (NLU) tasks?” To understand this, we conduct a study comparing the language representations trained using only the text to the language representations trained using vision supervision. More specifically, we compare the performance of the text encoders used in vision-language models to the vanilla pre-trained text encoders.

Few works [7, 27] evaluated the performance of vision-language and vanilla language models on GLUE. However, there exists a data discrepancy as these models are pre-trained on different domains of data making the comparisons unfair. To overcome this, we pre-train all the vanilla language models with the text captions used in multi-modal pre-training while keeping the identical training setting. Therefore, the only difference in training between vision-language and vanilla language models is the use of vision data.

For our experiments we use a diverse set of image-text models: ALBEF [12], BLIP [11] and METER [4] and video-text models: ALPRO [10], Frozen-in-time (FiT) [1] and VIOLET [5]. We evaluate these models on NLU benchmarks GLUE [29], SuperGlue [28] and Common sense reasoning datasets such as SocialIQA [25], CosmosQA [6], WinoGrande [23], CODAH [2] and HelLaSwag [33].

Our experiments show that (i) vision supervised language representations under perform compared to vanilla language representations on most of the Natural Language Understanding tasks like Natural Language Inference (NLI), sentence similarity, reading comprehension, linguistic probe and textual entailment. (ii) A similar trend is observed for commonsense reasoning benchmarks.

2. Related Work

Over the recent years there has been a tremendous progress in training vision and language together using large-scale multi-modal data. [3, 13, 14]. These models combine both the modalities into a single input and are trained using objectives similar to masked language modelling. Another line of work [1, 11, 12, 22] explore dual stream architectures in which there is a separate encoder for each of the modalities and the final representations are minimized using contrastive loss.

Natural Language Understanding involves several tasks such as text classification [19, 30], sentence similarity [20], Natural Language Inference [32] etc. However to evaluate the capability of models towards a broad range of NLU

tasks, benchmarks such as GLUE [29], Superglue [28] are introduced. Since then, these benchmarks are being used to comprehensively evaluate the performance of language models.

3. Experiments

3.1. Models

We experiment with a diverse set of image-text and video-text models. These models differ in the type of pre-training data used, in the architecture of the text encoder and in the sizes the text encoder. The comparison among the models is shown in the table 1.

3.1.1 ALBEF

ALBEF [12] is an image-text model pretrained on conceptual captions 12M (CC12M) [26], COCO [15], SBU captions [21] and visual genome [9]. It's text encoder has a pre-trained BERT [8] architecture with six transformer encoder layers.

3.1.2 BLIP

BLIP [10] is proposed as an extension to ALBEF model pretrained using the same data albeit with a large text encoder. It's text encoder has the same configuration as pre-trained BERT.

3.1.3 METER

METER [4] is an image-text model pretrained on conceptual captions 3M (CC3M), SBU captions and visual genome. Pre-trained RoBERTa [16] with six transformer encoder layers is used as the text encoder.

3.1.4 ALPRO

ALPRO [12] is a video-text model whose text encoder has a pre-trained BERT architecture with six transformer encoders. It is pre-trained on a combined data of conceptual captions 3M (CC3M) and WebVid-2M [1].

3.1.5 Frozen-in-time (FiT)

Frozen-in-time [1] is a dual stream transformer model pretrained on both image data conceptual captions 3M (CC3M) and video data WebVid-2M. DistilBERT [24] is used as the text encoder.

3.1.6 VIOLET

VIOLET [5] is a multi-modal transformer model pretrained end-to-end on YouTube 180M (YT180M) [34],

conceptual captions 3M (CC3M) and WebVid-2M. The text encoder follows the BERT architecture.

3.2. Datasets

For our analysis, we use GLUE, Superglue and common-sense reasoning datasets such as SocialQA, CosmosQA, WinoGrande, CODAH and HellaSwag. For all these datasets, we evaluate the models on the dev data.

3.3. Implementation

For fair comparison between the vision supervised text models and vanilla text models, we pre-train the vanilla text models with the text captions from the datasets used for large scale training of image-text and video-text models. Now, the only difference between these models is the use of vision data. We pre-train vanilla text models in the exact setup as the original vision-language models. We then fine-tune both the vision supervised text models and vanilla text models on downstream tasks. For GLUE, the maximum sentence length used is 200 and the models are trained for 5 epochs. In case of superglue, 250 is the maximum sentence length and the model are trained for 25 epochs. For commonsense reasoning, the models are trained for 10 epochs and 300 is the maximum sentence length. Unless otherwise stated, the results reported are the average of 5 runs.

4. Results

Table 2 shows the results on GLUE benchmark. From the tables, it is evident that vanilla language representations show superior performance compared to vision supervised language representations on most of the tasks across all the models. The drop in performance is significant for NLI tasks like MNLI and MNLI-mismatched (MNLI-mis). A similar trend is observed for sentence similarity (QQP), sentiment classification (SST2), reading comprehension (MRPC), linguistic probe (CoLA) and textual entailment (RTE). However, we see a huge improvement in performance for the Winograd NLI (WNLI) task.

Table 3 illustrates the results on superglue benchmark. From the table, we observe that vision supervised language representations under perform compared to vanilla language representations. For the tasks question answering (BoolQ), word in context (WiC), discourse (CB) we see a huge drop in performance. However, we see a significant improvement in performance for the casual reasoning (COPA) task. It is worth-noting that the performance is same for both the vanilla and vision supervised language representations on winograd schema challenge (WSC).

Table 4 demonstrates the results on commonsense reasoning datasets. As shown in the table, the performance of vanilla language representations surpass vision supervised language representations. There is a notable difference in performance on SocialQA, CosmosQA, WinoGrande and

Table 1. Comparison among different image-text and video-text models in-terms of pre-training data, architecture of the text encoders and size of the text encoder. CC denotes Conceptual captions [1], SBU denotes SBU captions [21] and VG represents visual genome [9].

Type	Model	Pre-training Data	Text Encoder	Num. layers
Image-text	ALBEF	CC12M + COCO + SBU + VG (14M)	BERT	6
	BLIP	CC12M + COCO + SBU + VG (14M)	BERT	12
	METER	CC3M + SBU + VG (4M)	RoBERTa	6
Video-text	ALPRO	CC3M + WebVid-2M (5M)	BERT	6
	FiT	CC3M + WebVid-2M (5M)	DistilBERT	6
	VIOLET	YT180M + CC3M + WebVid-2M (11M)	BERT	12

Table 2. Results on GLUE benchmark. MNLI-mis refers to the task MNLI mismatched and WNLI denotes the Winograd Schema Challenge. We see that language representations learnt through vision supervision under performs compared to vanilla language representations on all the tasks except WNLI.

Model	Type	MNLI	MNLI-mis	QQP	SST2	MRPC	CoLA	RTE	WNLI
ALBEF	Text	82.77	82.68	90.54	91.44	72.81	81.50	58.12	46.01
	Image-text	61.38	61.68	79.02	80.39	66.49	69.13	50.30	56.34
BLIP	Text	83.04	82.70	90.54	91.44	72.81	81.50	58.12	46.01
	Image-text	61.38	61.68	79.02	80.39	66.49	69.13	50.30	56.34
METER	Text	86.59	86.15	90.99	93.27	76.06	82.58	64.02	56.34
	Image-text	31.82	31.82	77.91	81.12	66.49	69.13	47.29	56.34
ALPRO	Text	82.96	82.81	90.64	92.05	70.96	79.93	60.41	45.07
	Video-text	62.53	63.26	79.35	80.96	66.49	69.13	54.39	56.34
FiT	Text	79.10	80.23	89.51	52.03	72.58	69.13	57.28	48.83
	Video-text	59.54	59.45	79.01	52.18	66.78	69.13	48.01	56.34
VIOLET	Text	83.19	83.59	90.68	92.74	71.92	81.66	59.93	52.58
	Video-text	61.38	61.68	79.02	80.39	66.49	69.13	50.30	56.34

Table 3. Results on Superglue benchmark. WiC represents Word-in-Context, CB represents CommitmentBank, COPA denotes Choice of Plausible Alternatives and WSC means The Winograd Schema Challenge.

Model	Type	BoolQ	WiC	CB	COPA	WSC
ALBEF	Text	70.41	63.13	76.79	48.00	63.46
	Image-text	63.30	55.02	63.93	51.60	63.46
BLIP	Text	70.41	63.13	76.43	48.00	63.46
	Image-text	63.30	55.02	63.93	51.60	63.46
METER	Text	72.40	66.11	75.00	46.80	63.46
	Image-text	66.87	53.98	69.64	50.80	63.46
ALPRO	Text	71.16	67.18	76.79	42.20	63.46
	Video-text	65.17	53.17	62.50	50.60	62.50
FiT	Text	68.91	62.38	69.29	44.80	63.46
	Video-text	64.69	53.20	70.71	53.80	63.46
VIOLET	Text	63.85	57.37	66.07	56.00	63.46
	Video-text	63.44	54.11	63.93	52.60	63.46

Table 4. Results on Commonsense reasoning tasks.

Model	Type	SocialQA	CosmosQA	WinoGrande	CODAH	HellaSwag
ALBEF	Text	40.50	26.45	53.12	25.72	25.04
	Image-text	33.47	25.24	49.57	25.72	24.48
BLIP	Text	52.27	25.72	56.88	26.02	25.24
	Image-text	33.47	25.24	49.57	25.72	24.48
METER	Text	58.39	31.32	59.59	24.40	25.04
	Image-text	33.47	25.00	49.57	25.72	24.48
ALPRO	Text	49.90	27.45	56.56	24.10	24.89
	Video-text	33.96	25.70	50.28	25.72	24.48
FiT	Text	45.46	30.87	56.75	25.12	26.54
	Video-text	33.35	25.77	50.33	24.52	24.59
VIOLET	Text	43.36	33.17	57.09	24.28	25.27
	Video-text	33.47	25.24	49.57	25.72	24.48

HellaSwag commonsense tasks. However for the CODA dataset, we observe vision supervised language representations outperform vanilla language representations for METER, ALPRO and VIOLET models.

5. CONCLUSION AND FUTURE DIRECTIONS

In this paper we comprehensively evaluated if the vision supervised language representations are beneficial to the language. We experimented with three image-text models ALBEF, BLIP, METER and three video-text models ALPRO, FiT, VIOLET on NLU benchmarks GLUE, superglue and commonsense reasoning tasks. Our experiments showed that vanilla language representations significantly outperform vision supervised language representations on most of the tasks. We believe these findings can shed light on the future directions to improve the vision-language pre-training that is beneficial to understanding the language.

References

- [1] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1728–1738, 2021.
- [2] Michael Chen, Mike D’Arcy, Alisa Liu, Jared Fernandez, and Doug Downey. Codah: An adversarially authored question-answer dataset for common sense. *arXiv preprint arXiv:1904.04365*, 2019.
- [3] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning.
- [4] Zi-Yi Dou, Yichong Xu, Zhe Gan, Jianfeng Wang, Shuohang Wang, Lijuan Wang, Chenguang Zhu, Pengchuan Zhang, Lu Yuan, Nanyun Peng, et al. An empirical study of training end-to-end vision-and-language transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18166–18176, 2022.
- [5] Tsu-Jui Fu, Linjie Li, Zhe Gan, Kevin Lin, William Yang Wang, Lijuan Wang, and Zicheng Liu. Violet: End-to-end video-language transformers with masked visual-token modeling. 2021.
- [6] Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Cosmos qa: Machine reading comprehension with contextual commonsense reasoning. 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 2391–2401, 2019.
- [7] Taichi Iki and Akiko Aizawa. Effect of visual extensions on natural language understanding in vision-and-language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2189–2196, 2021.
- [8] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186, 2019.
- [9] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73, 2017.
- [10] Dongxu Li, Junnan Li, Hongdong Li, Juan Carlos Niebles, and Steven CH Hoi. Align and prompt: Video-and-language pre-training with entity prompts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4953–4963, 2022.
- [11] Junnan Li, Dongxu Li, Caiming Xiong, and Steven CH Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. 2022.

- [12] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34:9694–9705, 2021.
- [13] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- [14] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXX 16*, pages 121–137. Springer, 2020.
- [15] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.
- [16] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [17] Avinash Madasu, Estelle Aflalo, Gabriela Ben Melech Stan, Shao-Yen Tseng, Gedas Bertasius, and Vasudev Lal. Improving video retrieval using multilingual knowledge transfer. In *Advances in Information Retrieval: 45th European Conference on Information Retrieval, ECIR 2023, Dublin, Ireland, April 2–6, 2023, Proceedings, Part I*, pages 669–684. Springer, 2023.
- [18] Avinash Madasu, Junier Oliva, and Gedas Bertasius. Learning to retrieve videos by asking questions. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 356–365, 2022.
- [19] Avinash Madasu and Vijjini Anvesh Rao. Sequential learning of convolutional features for effective text classification. 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 5658–5667, 2019.
- [20] Jonas Mueller and Aditya Thyagarajan. Siamese recurrent architectures for learning sentence similarity. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.
- [21] Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million captioned photographs. *Advances in neural information processing systems*, 24, 2011.
- [22] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.
- [23] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- [24] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- [25] Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social iqa: Commonsense reasoning about social interactions. 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 4463–4473, 2019.
- [26] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018.
- [27] Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15638–15650, 2022.
- [28] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems*, 32, 2019.
- [29] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *International Conference on Learning Representations*, 2018.
- [30] Sida I Wang and Christopher D Manning. Baselines and bigrams: Simple, good sentiment and topic classification. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 90–94, 2012.
- [31] Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. Simvlm: Simple visual language model pretraining with weak supervision. In *International Conference on Learning Representations*, 2021.
- [32] Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sentence understanding through inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122, 2018.
- [33] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, 2019.
- [34] Rowan Zellers, Ximing Lu, Jack Hessel, Youngjae Yu, Jae Sung Park, Jize Cao, Ali Farhadi, and Yejin Choi. Mer-

lot: Multimodal neural script knowledge models. *Advances in Neural Information Processing Systems*, 34:23634–23651, 2021.