

Hallucination Augmented Contrastive Learning for Multimodal Large Language Model

Chaoya Jiang¹ Haiyang Xu^{2*} Mengfan Dong¹ Jiaxing Chen¹ Wei Ye^{1*}
 Ming Yan² Qinghao Ye² Ji Zhang² Fei Huang² Shikun Zhang¹
¹National Engineering Research Center for Software Engineering, Peking University
²Alibaba Group

{jiangchaoya, wye}@pku.edu.cn, shuofeng.xhy@alibaba-inc.com

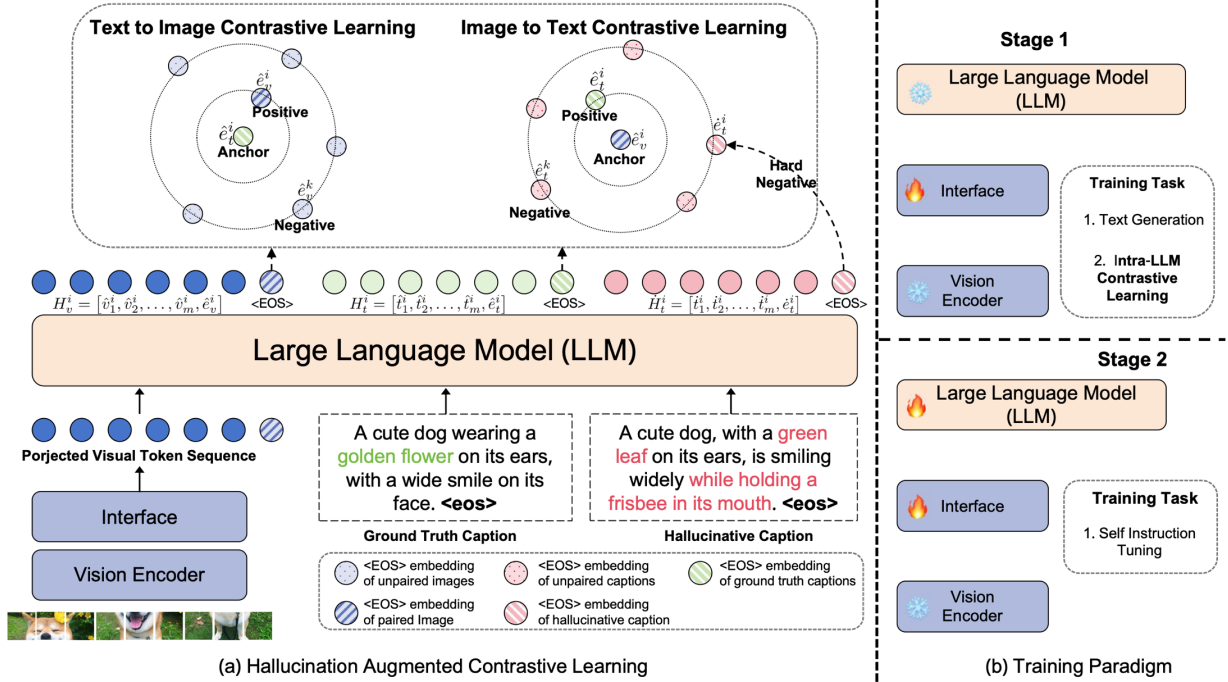


Figure 2. Subfigure (a) illustrates the proposed HACL. In this framework, we employ GPT-4 [38] to generate the hallucinative captions as the hard negative samples in the image-to-text contrastive learning. Subfigure (b) shows the training paradigm of HACL.

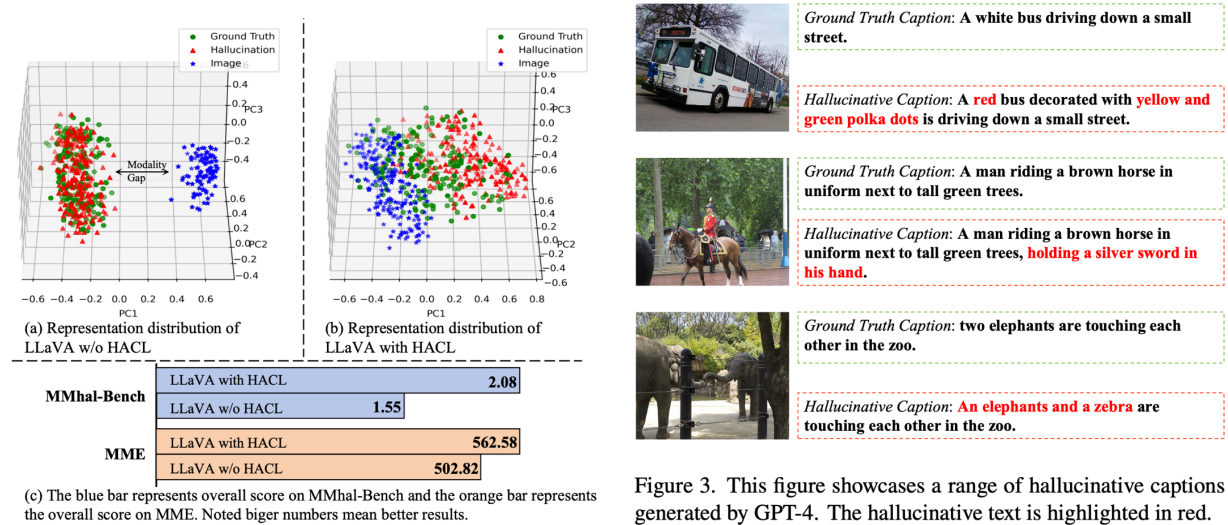


Figure 3. This figure showcases a range of hallucinative captions generated by GPT-4. The hallucinative text is highlighted in red.