## Energy-Guided Decoding for Object Hallucination Mitigation

Xixi Liu<sup>1</sup>, Ailin Deng<sup>2</sup>, Christopher Zach<sup>1</sup> <sup>1</sup> Chalmers University of Technology <sup>2</sup>National University of Singapore



CHALMERS **UNIVERSITY OF TECHNOLOGY** 



Unlike existing hallucination mitigation methods such as VCD [1] (which necessitates generating a sophisticated noisy version of the original visual inputs), OPERA [2](which relies on the beam-searching decoding mechanism), HALC [4] (requiring a pre-defined layer bucket and an external detector), and MMVP [3](relies on additional fine-tuning), our method is derived through the lens of internal states of a language decoder. It avoids the need of visual distortion, or fine-tuning visual encoders, or external detectors making it free from contrastive decoding. More importantly, the energy score at each layer can be computed with a single forward pass, making our method significantly less computationally demanding compared to OPERA [2] and HALC [4].

## Vision-Language Model Generation Selected Results The input tokens consisting of the visual tokens and language tokens is Our method consistently improves accuracy and F1 score on POPE benchmark denoted with x with the length of T. VLMs are commonly trained in with GQA dataset over three baseline methods with LLaVa-1.5. an autoregressive manner with a causal attention mask meaning that the prediction of the current token $x_t$ only depends on the previous tokens, LLaVA-1.5 InstructBLIP Datasets *Setting* formally, Yes ratio $\Delta_{gap} \downarrow$ Decodin Recall F1 Score↑ Recall F1 Score↑ Accuracy↑ Precision Accuracy↑ Precision Greedy 96.00 87.09 60.23 10.23 90.07 53.17 <u>3.17</u> 85.77 VCD [1] 93.40 83.34 62.07 12.07 80.90 77.35 87.40 82.07 56.50 6.50 81.33 75.24 54.37 HALC [3 85.90 96.00 87.19 60.10 10.10 85.97 83.08 90.33 86.55 4.37 $\mathbf{h} = \mathsf{VLM}(\mathbf{x}) = \{h_0, h_1, \cdots, h_{T-1}\},\$ (1)OPERA [2] 92.93 54.37 4.3787.33 87.23 87.47 **87.35** 50.13 0.13 89.05

where  $\mathbf{h}$  is the output state of the final layer of language decoder, and

the size of  $h_t$  is  $f_{dim}$ . The next token predictive distribution is defined as

(2) $p(x_t | x_{< t}) = \mathsf{Softmax}[\mathcal{H}(h_t)],$ 

where  $x_{<t}$  denotes the sequence of tokens before t-th position  $\{x_i\}_{i=0}^{t-1}$ and  $\mathcal{H} \in R^{f_{\text{dim}} \times V_{\text{size}}}$  is the learned vocabualry head.

## Methods

We use this energy score to identify the layer whose hidden state provides the most reliable representation of the input. In detail, the energy score is given by

$$\mathbf{Energy}(h_t^k) = -\operatorname{LogSumExp}[\mathcal{H}(h_t^k)]$$
(3)

where  $\mathcal{H}(h_t^k)$  denote the logits calculated at layer k for predicting token t. The layer  $k^* = \arg\min_k \mathbf{Energy}(h_t^k)$  with the lowest score is consequently selected for decoding.







## References

[1] Mitigating object hal- 603 lucinations in large vision-language models through visual 604 contrastive decoding. CVPR, 2024. [2] Opera: Alleviating hallucination in multi- 588 modal large language models via over-trust penalty and 589 retrospection-allocation. CVPR, 2024

[3] Eyes wide shut? exploring the661 visual shortcomings of multimodal llms. CVPR, 2024. [4] Halc: Object hallucination reduc-522 tion via adaptive focal-contrast decoding. ICML, 2024.

89.37

74.73

90.70

67.35

Energy (Ours)

87.73 **89.19** 

79.16

96.00

1.63

21.27

48.37

71.27

86.53

76.37

92.35

70.70

79.67 85.54

90.07 79.21

43.13

63.70

6.87

13.70

