

Bridging Compressed Image Latents and Multimodal Large Language Models

Chia-Hao Kao^{1,2}, Cheng Chien², Yu-Jen Tseng², Yi-Hsin Chen², Alessandro Gnutti¹, Shao-Yuan Lo³, Wen-Hsiao Peng², Riccardo Leonardi¹

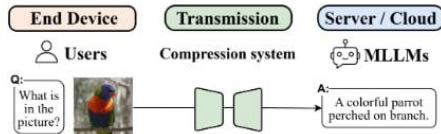
¹University of Brescia, Italy ²National Yang Ming Chiao Tung University, Taiwan ³Honda Research Institute, USA

Motivation

MLLMs excel at reasoning with text and images but typically require server-side hosting due to their size, necessitating **data transmission**

Compression is necessary:

w/o compression \rightarrow 24 bits per pixel (22MB for a 720p image)
with compression \rightarrow \approx 0.2 bits per pixel (\approx 0.15MB) \leftarrow **120x smaller**



Codecs optimized for human vision hinders MLLM performance



Existing coding for machines methods do not consider MLLMs

Backpropagation through task network renders **infeasible training**



Call for an efficient compression system designed for MLLMs

Contributions

This marks the **first exploration** of neural image coding for MLLMs

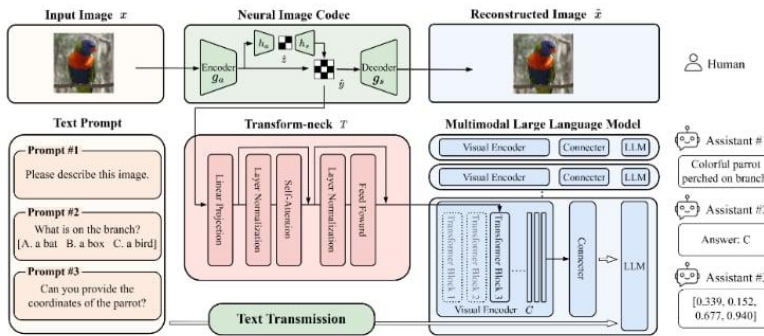
It adapts image latents directly to MLLMs **saving computational complexity**, while avoiding backpropagating through the heavy MLLM

It is broadly applicable to a wide range of codecs, tasks, and MLLMs

It is able to accommodate **various application scenarios** that involve human perception, machine perception, or both

Architecture

Directly **adapt compressed image latent representations** to connect with MLLMs

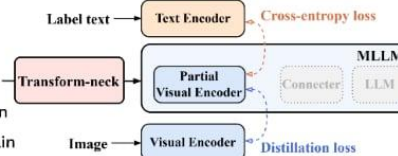


Surrogate Loss

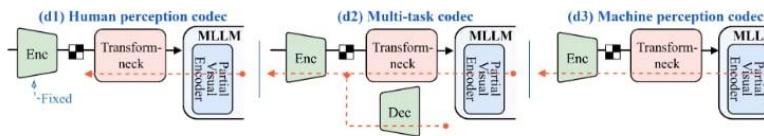
Leverage **only visual encoder**, avoiding backpropagation through entire MLLM

Distillation loss ensures alignment of the visual features before and after compression

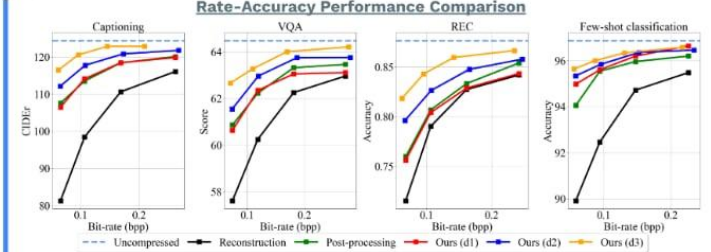
CE loss bridges features with the text domain



Application Scenarios



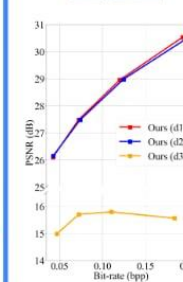
Experimental Results



Computational Complexity Comparison

Method	Component	Params (M)	kMAC/pixel
Ours (d1, d2, or d3)	Transform-neck	13.19	52.795
Post-processing	Decoder	7.34	112.00
	Post-processing network	31.04	835.72
	First 2 layers of visual encoder	25.78	70.24

Rate-Distortion Comparison



Visualization

