

VCR: A Task for Pixel-Level Complex Reasoning in Vision Language Models via Restoring Occluded Text

Tianyu Zhang*, Suyuchen Wang*, Lu Li, Ge Zhang,
Perouz Taslakian, Sai Rajeswar, Jie Fu, Bang Liu, Yoshua Bengio

ICLR 2025

Visual Caption Restoration (VCR) Overview

In this paper, we introduce:

1. A Task: Visual Caption Restoration (VCR)
2. A Dataset + Benchmark: VCR-Wiki
3. A Hidden Test: VCR-Hidden

“What are the covered texts in the image?”



Lincoln is the luxury vehicle division of American automobile manufacturer Ford. Marketed among the top luxury vehicle brands in the United States, for

The VCR Task

Easy



Lincoln is the luxury vehicle division of American automobile manufacturer Ford. Marketed among the top luxury vehicle brands in the United States, for

Hard



Lincoln is the luxury vehicle division of American automobile manufacturer Ford. Marketed among the top luxury vehicle brands in the United States, for

English



Chinese

科西嘉是地中海西部的一座岛屿，也属于法国领土集体，位于意大利半岛的西南面，最近的地块距离意大利半岛约200公里。科西嘉岛三分之二



科西嘉是地中海西部的一座岛屿，也属于法国领土集体，位于意大利半岛的西南面，最近的地块距离意大利半岛约200公里。科西嘉岛三分之二

The VCR-Wiki Dataset

Latest data source, multi-lingual, multi-fonts, multi-resolution, multi-layout

Should people with bigger cars pay more for parking? As part of a public consultation, it is asking for views on whether residents with larger vehicles should pay more for permits.

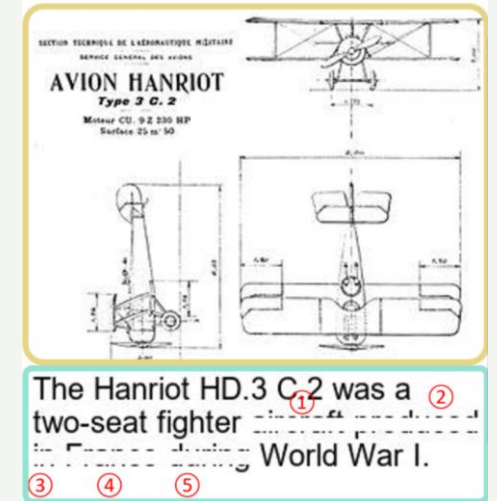


The VCR-Hidden Test

The Visual Caption Restoration (VCR) Task

1. A task: Visual Caption Restoration (VCR)

- Task: given an image with two parts: a visual element and partly-occluded text-in-image, restore the covered text
- A typical thinking process for this task involves using pixel-level hints for texts and multi-step reasoning



• Characteristics:

- Unsolvable by OCR: the remaining pixels of the covered texts are crucial and hard to use
- Have deterministic answer: the remaining pixels gives definitive answer; No LLM-as-a-judge needed
- Flexibility: over covered area size, text to cover...

1. This image is **very likely** about a manual of an old fashion plane.
2. The visible pixel level hint of ① shows that it is a word ending with "ft". Given the previous word is "fighter", the first blank is likely "aircraft".
3. The title of the manual mentioned "HANRIOT" which is a French manufacturer.
4. ④ start from "F" thus, combined with 4, we induce it should be "France".
5. It is very likely that we need a "adj." between "France" and "World War I". Given ⑤ ends with "g" and start with "d". We can infer it is "during".
6. ② start from "p" and end with "d". It needs to be a verb. Connecting with 4, it is "produced".
7. ③ is a word starting from "i" and only has 2 character. It evident that it is "in".

The VCR-Wiki Dataset

2. A Dataset + Benchmark: VCR-Wiki
- Sourced from Wikipedia dump
 - Multilingual: Contains En and Zh splits
 - Multi-difficulty: Contains easy and hard splits:
 - Easy: The text cover causes OCR to completely fail
 - Hard: Only 1-2 pixels left on the top and bottom
 - The easy / hard performance are directly comparable, with exact same splits and texts and only pixel-level difference on the cover
 - Contains train / val / test sets
 - SFT on the training set can improve real-world occluded text recognition

Easy



Lincoln is the luxury vehicle division of American automobile manufacturer Ford. Marketed among the top luxury vehicle brands in the United States, for

Hard



Lincoln is the luxury vehicle division of American automobile manufacturer Ford. Marketed among the top luxury vehicle brands in the United States, for



科西嘉是地中海西部的一座岛屿，也是法国的一个领土集体，位于法国大陆部分的东南面，意大利半岛的西面，最近的地块是紧邻意大利的撒丁岛。科西嘉岛三分之二



科西嘉是地中海西部的一座岛屿，也是法国的一个领土集体，位于法国大陆部分的东南面，意大利半岛的西面，最近的地块是紧邻意大利的撒丁岛。科西嘉岛三分之二

The VCR-Hidden Hidden Test

3. A Hidden Test: VCR-Hidden
 - To avoid data contamination on VCR-Wiki's test set, we created a hidden test for VCR:
 - Up-to-date source: we use latest news to update the data source
 - More flexible data: the data contains multi-lingual, multi-fonts, multi-resolution, multi-layout VCR images for maximum flexibility

Should people with higher car payments for parking? As part of a public consultation, it is asking for views on whether residents with higher car payments should pay more for permits.



American star Mikaela Shiffrin suffered a ~~penetration~~ ^{penetration} wound to the right side of her abdomen and severe muscle trauma. The two-time Olympic champion had been seeking a record-extending 100th World Cup win and looked well placed to reach the milestone before ~~falling and sustaining a knee injury~~ ^{falling and sustaining a knee injury}, netting on her second run.



Huge Gap Between SoTA VLM and Human

Table 4: Human evaluation results on the VCR task for in terms of exact matches. N is the number of puzzles in each language.

	EN Easy (N = 169)		EN Hard (N = 169)		ZH Easy (N = 188)		ZH Hard (N = 188)	
	Mean (%)	SD (%)	Mean (%)	SD (%)	Mean (%)	SD (%)	Mean (%)	SD (%)
All	96.65	0.34	91.12	1.18	98.58	0.31	91.84	0.81
Filtered	98.62	0.34	97.63	2.13	99.47	0.00	96.63	1.11

Native speakers can very easily achieve
~100% restoration accuracy on VCR-Wiki!

Language	Mode	Model name	Model size	Exact match (%) ↑			Jaccard index (%) ↑		
				VI + TEI	TEI	Δ	VI + TEI	TEI	Δ
English	Hard	Closed-source models							
		Claude 3 Opus	-	37.80.28	50.00.33	-12.2	57.680.8	70.160.64	-12.48
		Claude 3.5 Sonnet	-	41.741.69	44.721.78	-2.98	56.151.46	58.541.6	-2.4
		Gemini 1.5 Pro	-	28.071.58	38.761.68	-10.68	51.912.22	59.621.27	-7.72
		GPT-4 Turbo	-	45.150.28	48.640.57	-3.5	65.720.25	67.860.2	-2.14
		GPT-4o	-	73.20.16	82.430.17	-9.22	86.170.21	92.010.2	-5.84
		GPT-4V	-	25.830.44	14.950.3	10.87	44.630.48	30.080.67	14.56
		Qwen-VL-Max	-	41.650.32	52.720.2	-11.07	61.180.35	70.190.37	-9.01
		Reka Core	-	6.710.89	11.181.15	-4.47	25.840.95	35.831.05	-9.99
		Open-source models							
		CogVLM2	19B	37.980.18	17.680.06	20.3	59.990.05	39.690.03	20.3
		DeepSeek-VL	1.3B	0.160.01	0.390.02	-0.23	11.890.02	11.470.03	0.42
		DeepSeek-VL	7B	1.00.02	1.750.03	-0.75	15.90.08	17.20.04	-1.3
		DocOwl-1.5-Omni	8B	0.040.0	0.020.0	0.01	7.760.01	7.740.02	0.03
		Monkey	7B	1.960.04	2.430.03	-0.48	14.020.03	14.110.03	-0.09
		Idolies2	8B	0.650.01	0.940.02	-0.29	9.930.05	12.570.02	-2.64
		InternLM-XComposer2-VL	7B	0.70.01	0.920.01	-0.22	12.510.02	13.230.02	-0.72
		InternVL-V1.5	25.5B	1.990.02	6.490.04	-4.5	16.730.06	26.40.03	-9.67
		MiniCPM-V2.5	8B	1.410.03	1.960.02	-0.55	11.940.02	13.370.04	-1.43
		Qwen-VL	7B	2.00.03	2.320.03	-0.32	15.040.05	14.270.05	0.77
		Yi-VL	34B	0.070.0	0.050.0	0.02	4.310.02	5.890.02	-1.58

English
Hard

Chinese	Hard	Closed-source models							
		Claude 3 Opus	-	0.3 _{0.18}	0.1 _{0.1}	0.2	9.22 _{0.38}	8.09 _{0.33}	1.13
		Claude 3.5 Sonnet	-	0.2 _{0.15}	0.0 _{0.0}	0.2	4.0 _{0.33}	2.37 _{0.23}	1.63
		Gemini 1.5 Pro	-	0.7 _{0.26}	0.5 _{0.23}	0.2	11.82 _{0.51}	11.75 _{0.44}	0.07
		GPT-4o	-	2.2 _{0.47}	1.8 _{0.4}	0.4	22.7 _{20.67}	22.8 _{90.05}	-0.17
		GPT-4 Turbo	-	0.0 _{0.0}	0.2 _{0.13}	-0.2	8.5 _{80.3}	6.87 _{0.28}	1.72
		Qwen-VL-Max	-	0.89 _{0.06}	1.38 _{0.1}	-0.49	5.4 _{0.19}	12.2 _{0.18}	-6.89
		Reka Core	-	0.0 _{0.0}	0.0 _{0.0}	0	3.35 _{0.23}	2.97 _{0.2}	0.38
		Open-source models							
		CogVLM2-Chinese	19B	1.34 _{0.03}	2.67 _{0.02}	-1.32	17.35 _{0.03}	19.51 _{0.03}	-2.16
		CogVLM2-Chinese-FT	19B	42.11 _{0.09}	45.63 _{0.06}	-3.51	65.67 _{0.15}	69.28 _{0.04}	-3.61
		DeepSeek-VL	1.3B	0.0 _{0.0}	0.0 _{0.0}	0	6.46 _{0.01}	3.22 _{0.02}	3.24
		DeepSeek-VL	7B	0.0 _{0.0}	0.0 _{0.0}	0	5.11 _{0.01}	7.21 _{0.01}	-2.1
		DocOwl-1.5-Omni	8B	0.0 _{0.0}	0.0 _{0.0}	0	1.37 _{0.01}	4.07 _{0.02}	-2.7
		Monkey	7B	0.12 _{0.01}	0.07 _{0.0}	0.05	6.36 _{0.01}	6.68 _{0.03}	-0.32
		InternLM-XComposer2-VL	7B	0.07 _{0.01}	0.09 _{0.0}	-0.02	8.97 _{0.02}	8.51 _{0.01}	0.46
		InternVL-V1.5	25.5B	0.03 _{0.0}	0.1 _{0.01}	-0.07	8.46 _{0.01}	6.27 _{0.04}	2.19
		MiniCPM-V2.5	8B	0.09 _{0.0}	0.08 _{0.0}	0.01	7.39 _{0.02}	7.89 _{0.01}	-0.5
		MiniCPM-V2.5-FT	8B	1.53 _{0.01}	1.11 _{0.02}	0.42	18.0 _{0.03}	15.35 _{0.02}	2.65
		Qwen-VL	7B	0.01 _{0.0}	0.01 _{0.0}	0	1.17 _{0.01}	0.12 _{0.0}	1.06
		Yi-VL	34B	0.0 _{0.0}	0.0 _{0.0}	0	4.12 _{0.0}	1.81 _{0.01}	2.31
		Yi-VL	6B	0.0 _{0.0}	0.0 _{0.0}	0	4.0 _{0.01}	1.88 _{0.01}	2.12

Chinese
Hard

... Yet a lot of SoTA VLMs achieve
nearly 0 restoration accuracy on VCR-Wiki

More interesting findings
(and VLM development guidelines) in the paper 

See You in Poster Session 2!



arXiv



GitHub



Hugging Face



OpenReview



Test with EvolvingLMMs-Lab/Imms-eval

[EvolvingLMMs-Lab/Imms-eval:](#)
[Accelerating the development of large multimodal models \(LMMs\) with Imms-eval \(github.com\)](#)



Test with open-compass/VLM EvalKit

[open-compass/VLM EvalKit:](#) Open-source evaluation toolkit of large vision-language models (LVLMs), support GPT-4v, Gemini, QwenVLPlus, 50+ HF models, 20+ benchmarks (github.com)