



EViL-Probe

A Composite Benchmark for Extensive Visio-Linguistic Probing

Marie Bexte¹ | Andrea Horbach^{1,2} | Torsten Zesch¹ ¹FernUniversität in Hagen, Germany | ²Hildesheim University, Germany





CATALPA – Center of Advanced Technology for Assisted Learning and Predictive Analytics



How well can visio-linguistic models align descriptions with photographs ?

Systematically manipulate target words in image descriptions



CATALPA



Visio-Linguistic Probing







EViL-Probe





CATALPA



CATALPA











Evaluation: Paired Accuracy

1 probe = 2 descriptions of the same image (1 match, 1 mismatch) processed as 2 separate classifications







Evaluation: Paired Accuracy





Evaluation: Paired Accuracy

Not anchored on another text

Respects text correspondence

Works with positive probes





	Lхмеrт	UNITER	ΛΙΓΓΑ	S оно	ALBEF	TcL	BLIP						
Negative Probes													
- Acc	.54	.62	.61	.57	.63	.62	.63						
- Acc _{paired}	.27	.26	.24	.20	.29	.27	.29						
Positive Probes													
- Acc	.47	.80	.83	.82	.54	.52	.75						
- Acc _{paired}	.25	.70	.75	.73	.41	.37	.65						

















	Lхмеrт	UNITER	VILLA	Ѕоно	ALBEF	TcL	ВLIP						
Negative Probes													
- Acc	.54	.62	.61	.57	.63	.62	.63						
- Acc _{paired}	.27	.26	.24	.20	.29	.27	.29						
Positive Probes													
- Acc	.47	.80	.83	.82	.54	.52	.75						
- Acc _{paired}	.25	.70	.75	.73	.41	.37	.65						





	Lхмеrт	UNITER	VILLA	Ѕоно	ALBEF	TcL	ВыР						
Negative Probes													
- Acc	.54	.62	.61	.57	.63	.62	.63						
- Acc _{paired}	.27	.26	.24	.20	.29	.27	.29						
Positive Probes													
- Acc	.47	.80	.83	.82	.54	.52	.75						
- Acc _{paired}	.25	.70	.75	.73	.41	.37	.65						







	- LXMERT	- UNITER	- VILLA	OHOS -	- ALBEF	- TCL	- BLIP	- LXMERT	- UNITER	- VILLA	OHOS -	- ALBEF	- TCL	- BLIP	
Attribute -	.52	.61	.60	.55	.62	.60	.62	.26	.25	.23	.15	.27	.24	.26	
Color -	.55	.66	.65	.58	.68	.66	.72	.30	.34	.32	.23	.38	.35	.46	
Image Type -	.52	.56	.55	.51	.76	.66	.75	.23	.13	.10	.10	.52	.33	.52	
Negation -	.51	.53	.54	.52	.50	.51	.51	.23	.09	.09	.09	.03	.05	.08	
Noun -	.59	.76	.75	.67	.71	.70	.73	.32	.53	.52	.38	.45	.42	.49	
Number –	.51	.57	.56	.51	.52	.53	.52	.23	.17	.15	.09	.08	.09	.07	
Random -	.73	.95	.97	.94	.97	.96	.98	.47	.91	.94	.88	.95	.92	.96	
Semantic Role -	.50	.51	.50	.50	.52	.53	.52	.25	.07	.04	.06	.09	.13	.07	
Spatial Relation -	.51	.52	.51	.50	.51	.52	.50	.24	.06	.05	.04	.05	.08	.04	
Verb -	.52	.61	.60	.53	.62	.61	.62	.25	.26	.23	.11	.28	.26	.26	
Video-Based -	.50	.52	.52	.51	.54	.53	.53	.18	.15	.15	.14	.18	.16	.15	
Word Order -	.50	.53	.52	.50	.54	.53	.53	.22	.07	.06	.06	.12	.11	.09	
Hypernyms –	.42	.70	.75	.70	.37	.36	.60	.20	.56	.62	.55	.23	.23	.45	
Paraphrase -	.46	.78	.82	.80	.45	.44	.71	.24	.73	.76	.72	.38	.37	.65	
Perspective -	.48	.85	.87	.89	.69	.68	.88	.23	.71	.75	.80	.40	.38	.76	
Slang -	.42	.69	.71	.70	.30	.28	.62	.19	.52	.55	.57	.16	.14	.47	
Specificity -	.59	.99	.99	.94	.88	.81	.90	.37	.98	.98	.90	.77	.63	.81	
Word Order -	.45	.77	.82	.85	.54	.51	.76	.27	.70	.79	.79	.48	.45	.74	
	Acc								Acc_{paired}						





	- LXMERT	- UNITER	- VILLA	OHOS -	- ALBEF	- TCL	- BLIP	- LXMERT	- UNITER	- VILLA	OHOS -	- ALBEF	- TCL	- BLIP
Attribute -	.52	.61	.60	.55	.62	.60	.62	.26	.25	.23	.15	.27	.24	.26
Color -	.55	.66	.65	.58	.68	.66	.72	.30	.34	.32	.23	.38	.35	.46
Image Type -	.52	.56	.55	.51	.76	.66	.75	.23	.13	.10	.10	.52	.33	.52
Negation -	.51	.53	.54	.52	.50	.51	.51	.23	.09	.09	.09	.03	.05	.08
Noun -	.59	.76	.75	.67	.71	.70	.73	.32	.53	.52	.38	.45	.42	.49
Number -	.51	.57	.56	.51	.52	.53	.52	.23	.17	.15	.09	.08	.09	.07
Random -	.73	.95	.97	.94	.97	.96	.98	.47	.91	.94	.88	.95	.92	.96
Semantic Role -	.50	.51	.50	.50	.52	.53	.52	.25	.07	.04	.06	.09	.13	.07
Spatial Relation -	.51	.52	.51	.50	.51	.52	.50	.24	.06	.05	.04	.05	.08	.04
Verb -	.52	.61	.60	.53	.62	.61	.62	.25	.26	.23	.11	.28	.26	.26
Video-Based –	.50	.52	.52	.51	.54	.53	.53	.18	.15	.15	.14	.18	.16	.15
Word Order -	.50	.53	.52	.50	.54	.53	.53	.22	.07	.06	.06	.12	.11	.09
Hypernyms -	.42	.70	.75	.70	.37	.36	.60	.20	.56	.62	.55	.23	.23	.45
Paraphrase -	.46	.78	.82	.80	.45	.44	.71	.24	.73	.76	.72	.38	.37	.65
Perspective -	.48	.85	.87	.89	.69	.68	.88	.23	.71	.75	.80	.40	.38	.76
Slang -	.42	.69	.71	.70	.30	.28	.62	.19	.52	.55	.57	.16	.14	.47
Specificity -	.59	.99	.99	.94	.88	.81	.90	.37	.98	.98	.90	.77	.63	.81
Word Order -	.45	.77	.82	.85	.54	.51	.76	.27	.70	.79	.79	.48	.45	.74
				Acc				Acc_{paired}						





	- LXMERT	- UNITER	- VILLA	OHOS-	- ALBEF	- TCL	- BLIP	- LXMERT	- UNITER	- VILLA	OHOS -	- ALBEF	- TCL	- BLIP
Attribute -	.52	.61	.60	.55	.62	.60	.62	.26	.25	.23	.15	.27	.24	.26
Color -	.55	.66	.65	.58	.68	.66	.72	.30	.34	.32	.23	.38	.35	.46
Image Type -	.52	.56	.55	.51	.76	.66	.75	.23	.13	.10	.10	.52	.33	.52
Negation $-$.51	.53	.54	.52	.50	.51	.51	.23	.09	.09	.09	.03	.05	.08
Noun -	.59	.76	.75	.67	.71	.70	.73	.32	.53	.52	.38	.45	.42	.49
Number -	.51	.57	.56	.51	.52	.53	.52	.23	.17	.15	.09	.08	.09	.07
Random -	.73	.95	.97	.94	.97	.96	.98	.47	.91	.94	.88	.95	.92	.96
Semantic Role -	.50	.51	.50	.50	.52	.53	.52	.25	.07	.04	.06	.09	.13	.07
Spatial Relation -	.51	.52	.51	.50	.51	.52	.50	.24	.06	.05	.04	.05	.08	.04
Verb -	.52	.61	.60	.53	.62	.61	.62	.25	.26	.23	.11	.28	.26	.26
Video-Based –	.50	.52	.52	.51	.54	.53	.53	.18	.15	.15	.14	.18	.16	.15
Word Order -	.50	.53	.52	.50	.54	.53	.53	.22	.07	.06	.06	.12	.11	.09
Hypernyms –	.42	.70	.75	.70	.37	.36	.60	.20	.56	.62	.55	.23	.23	.45
Paraphrase -	.46	.78	.82	.80	.45	.44	.71	.24	.73	.76	.72	.38	.37	.65
Perspective -	.48	.85	.87	.89	.69	.68	.88	.23	.71	.75	.80	.40	.38	.76
Slang -	.42	.69	.71	.70	.30	.28	.62	.19	.52	.55	.57	.16	.14	.47
Specificity -	.59	.99	.99	.94	.88	.81	.90	.37	.98	.98	.90	.77	.63	.81
Word Order -	.45	.77	.82	.85	.54	.51	.76	.27	.70	.79	.79	.48	.45	.74
Acc								Acc_{paired}						





	- LXMERT	- UNITER	- VILLA	OHOS -	- ALBEF	- TCL	- BLIP	- LXMERT	- UNITER	- VILLA	OHOS -	- ALBEF	- TCL	- BLIP	
Attribute -	.52	.61	.60	.55	.62	.60	.62	.26	.25	.23	.15	.27	.24	.26	
Color -	.55	.66	.65	.58	.68	.66	.72	.30	.34	.32	.23	.38	.35	.46	
Image Type -	.52	.56	.55	.51	.76	.66	.75	.23	.13	.10	.10	.52	.33	.52	
Negation -	.51	.53	.54	.52	.50	.51	.51	.23	.09	.09	.09	.03	.05	.08	
Noun -	.59	.76	.75	.67	.71	.70	.73	.32	.53	.52	.38	.45	.42	.49	
Number –	.51	.57	.56	.51	.52	.53	.52	.23	.17	.15	.09	.08	.09	.07	
Random -	.73	.95	.97	.94	.97	.96	.98	.47	.91	.94	.88	.95	.92	.96	
Semantic Role -	.50	.51	.50	.50	.52	.53	.52	.25	.07	.04	.06	.09	.13	.07	
Spatial Relation -	.51	.52	.51	.50	.51	.52	.50	.24	.06	.05	.04	.05	.08	.04	
Verb -	.52	.61	.60	.53	.62	.61	.62	.25	.26	.23	.11	.28	.26	.26	
Video-Based –	.50	.52	.52	.51	.54	.53	.53	.18	.15	.15	.14	.18	.16	.15	
Word Order -	.50	.53	.52	.50	.54	.53	.53	.22	.07	.06	.06	.12	.11	.09	
Hypernyms –	.42	.70	.75	.70	.37	.36	.60	.20	.56	.62	.55	.23	.23	.45	
Paraphrase -	.46	.78	.82	.80	.45	.44	.71	.24	.73	.76	.72	.38	.37	.65	
Perspective -	.48	.85	.87	.89	.69	.68	.88	.23	.71	.75	.80	.40	.38	.76	
Slang -	.42	.69	.71	.70	.30	.28	.62	.19	.52	.55	.57	.16	.14	.47	
Specificity -	.59	.99	.99	.94	.88	.81	.90	.37	.98	.98	.90	.77	.63	.81	
Word Order -	.45	.77	.82	.85	.54	.51	.76	.27	.70	.79	.79	.48	.45	.74	
	Acc								Acc_{paired}						







Conclusions:

Benchmark is challenging for all tested models Positive probes help disambiguate performance

