

DUET: Cross-modal Semantic Grounding for Contrastive Zero-shot Learning

Zhuo Chen^{1, 5}, Yufeng Huang^{2, 5}, Jiaoyan Chen³, Yuxia Geng^{1, 2}, Wen Zhang^{2, 6}, Yin Fang^{1, 5}, Jeff Z. Pan⁴, Huajun Chen^{1, 5*}



1. College of Computer Science and Technology, Zhejiang University
2. School of Software Technology, Zhejiang University
3. Department of Computer Science, The University of Manchester
4. School of Informatics, The University of Edinburgh
5. Alibaba-Zhejiang University Joint Institute of Frontier Technologies

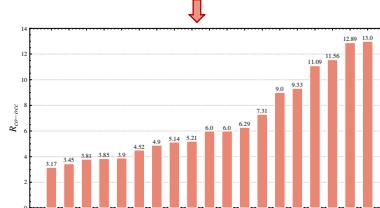
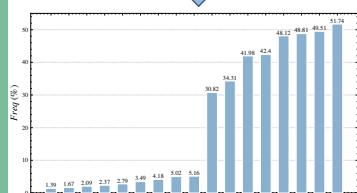


1 Challenges

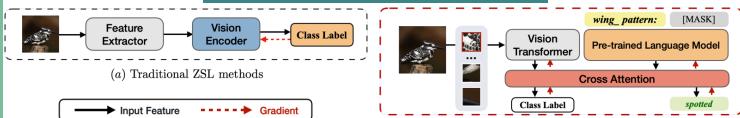


(a) Imbalance distribution of attributes

(b) Attribute co-occurrence

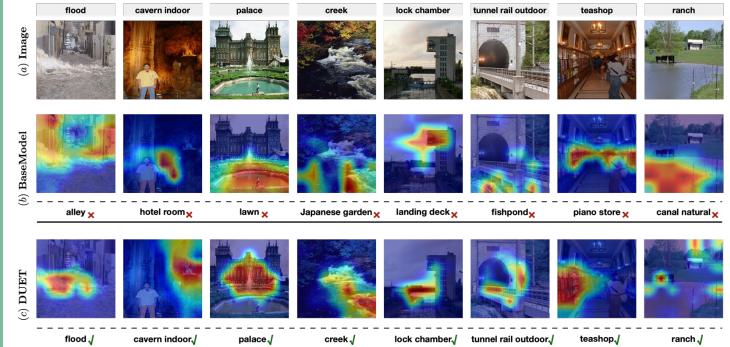


2 ZSL Paradigm comparison

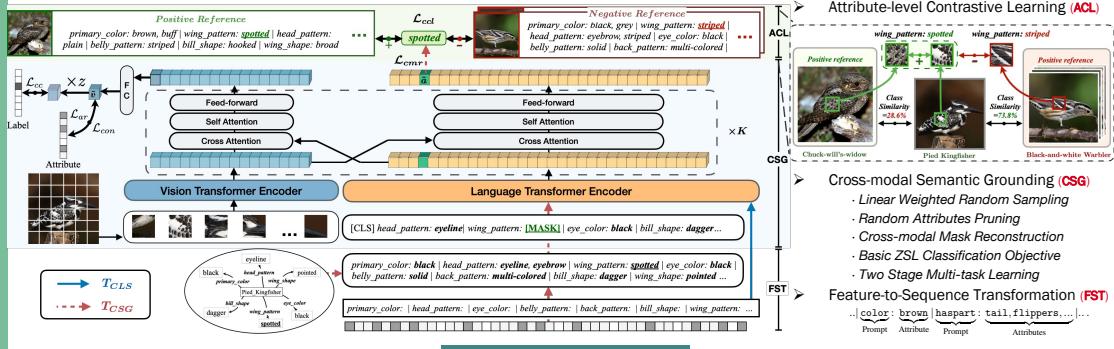


We exploits the semantics of PLMs to augment the transformer-based vision encoder via reconstructing masked attributes with a cross-model attention mechanism.

5 Semantic Grounding and Visualization



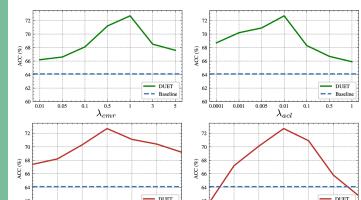
3 Propose DUET



4 Experiments

Image Class	Prompt	Attribute	CUB			SUN			AWA2					
			CZSL	GZSL	TI	CZSL	GZSL	TI	CZSL	GZSL	TI			
wing (S)	transportation	brown	64.9	52.8	64.7	58.1	66.0	45.6	40.7	43.0	72.2	59.8	75.1	66.6
coarse material	coarse	more (39.81%)	50.8	37.5	57.55%	50.8	44.1	55.1	22.0	31.4	71.5	62.1	77.3	68.8
coarse material	GT	osseous (42.21%)	69.4	56.4	63.8	59.9	62.6	55.1	22.0	31.4	71.5	62.1	77.3	68.8
specific material	fire (40.75%)	waves (27.78%)	63.1	66.8	65.3	63.3	48.8	38.6	43.1	70.4	63.1	78.6	70.0	
specific material	GT	osseous (42.21%)	67.5	63.1	69.7	60.3	47.9	37.8	42.2	70.4	63.1	78.6	70.0	
environment	function	nesting (45.42%)	61.0	59.7	60.3	59.7	62.2	50.3	37.8	42.2	60.4	75.1	67.0	
environment	GT	nesting (45.42%)	55.7	59.9	57.7	57.7	47.4	37.2	41.7	60.4	75.4	67.1		
function	FREE (ICCV) (2021)	climb (28.55%)	55.7	59.9	57.7	57.7	47.4	37.2	41.7	60.4	75.4	67.1		
function	HSVA (NeurIPS) (2021)	flowers (10.89%)	62.8	52.7	58.3	55.3	63.8	48.6	39.0	43.3	59.3	76.6	66.8	
function	AGZSL (ICLR) (2021)	scouting (9.37%)	57.2	41.4	49.7	45.2	63.3	39.9	40.2	34.3	73.8	65.1	78.9	71.3
basketball arena (S)	function	reports (48.15%)	72.0	65.3	69.3	67.2	61.6	41.9	34.0	37.6	68.4	57.1	72.4	63.9
basketball arena (S)	GT	congregating (47.13%)	53.2	50.7	59.8	57.1	45.0	37.2	40.7	67.7	63.6	70.8	67.0	
environment	function	competing (49.05%)	66.0	56.7	59.6	58.1	59.4	52.3	24.3	33.2	67.9	60.3	75.7	67.1
environment	GT	competing (49.05%)	77.5	63.1	73.5	66.1	63.8	44.0	31.7	36.8	73.6	67.1	76.5	71.5
surface	function	glossy (41.56%)	76.1	60.0	73.5	66.1	63.8	44.0	31.7	36.8	73.6	67.1	76.5	71.5
surface	GT	glossy (41.56%)	77.8	64.8	69.3	67.2	62.8	38.1	35.7	36.9	67.3	64.8	77.5	70.6
bus depot (S)	transportation	driving (42.32%)	76.1	68.7	67.5	68.1	65.8	52.2	34.2	41.3	70.1	62.0	74.5	67.7
transportation	GT	driving (42.32%)	76.8	69.3	68.3	68.8	65.6	52.6	33.4	40.8	70.1	61.3	82.3	70.2
material	GT	asphalt (57.30%)	67.6	50.9	55.0	52.0	51.9	41.3	20.5	24.3	67.3	59.3	76.6	66.8
material	GT	natural (54.51%)	76.8	69.3	68.3	68.8	65.6	52.6	33.4	40.8	70.1	61.3	82.3	70.2
light	GT	direct sunray (20.36%)	72.3	62.9	72.8	67.5	64.4	45.7	45.8	45.8	69.9	63.7	84.7	72.7

(a) Attribute prediction for interpretation



(c) Hyperparameter analysis of coefficients

(b) Overall results on standard ZSL datasets compared with 14 representative methods

