



# **Text-Enhanced Attribute-Based Attention for Generalized Zero-Shot Fine-Grained Image Classification**

Yan-He Chen and Mei-Chen Yeh

Department of Computer Science and Information Engineering  
National Taiwan Normal University

# Generalized zero-shot learning

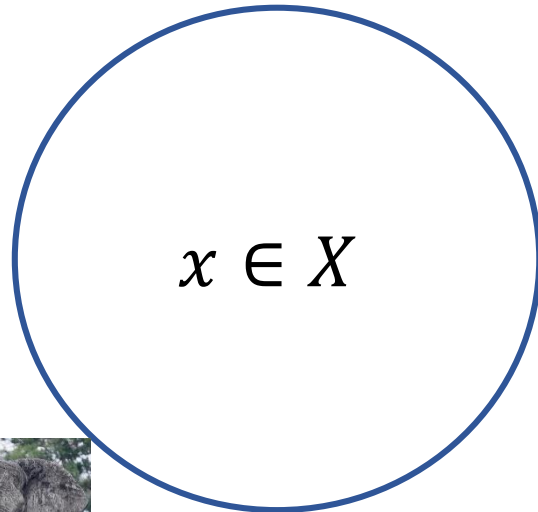
Training set (seen)



fox



elephant



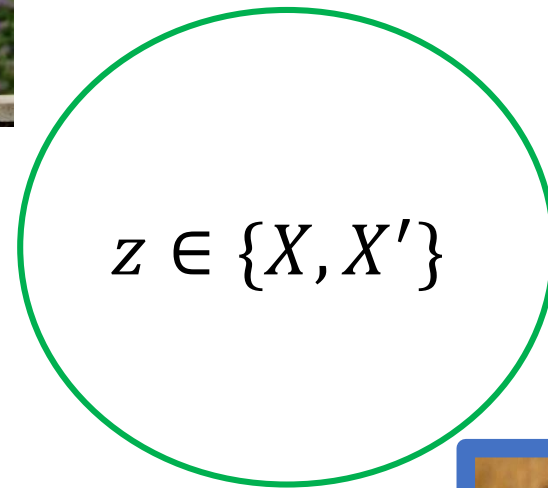
Test set (seen & unseen)



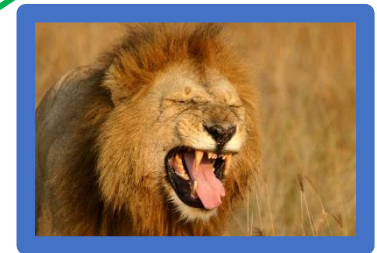
fox



horse



elephant



lion

$$X \cap X' = \emptyset$$

# Fine-grained image classification

Small inter-class variations

**Arctic Tern**



**Caspian Tern**



# Fine-grained image classification

Large intra-class variations



Laysan albatross: Adult



Laysan albatross: Juvenile

# Attributes for zero-shot learning

	Artic Tern	Caspian Tern
has_bill_shape::dagger	<input type="checkbox"/>	<input checked="" type="checkbox"/>
has_bill_shape::all-purpose	<input checked="" type="checkbox"/>	<input type="checkbox"/>
has_wing_color::grey	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
has_upperparts_color::white	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
has_leg_color::red	<input checked="" type="checkbox"/>	<input type="checkbox"/>
has_breast_pattern::solid	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
has_forehead_color::black	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
has_under_tail_color::white	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>



# Contributions of the paper

- For each class, a classification model should focus on the most relevant attributes.
- We propose to leverage auxiliary information in the form of *text descriptions* to achieve the goal.



Adult

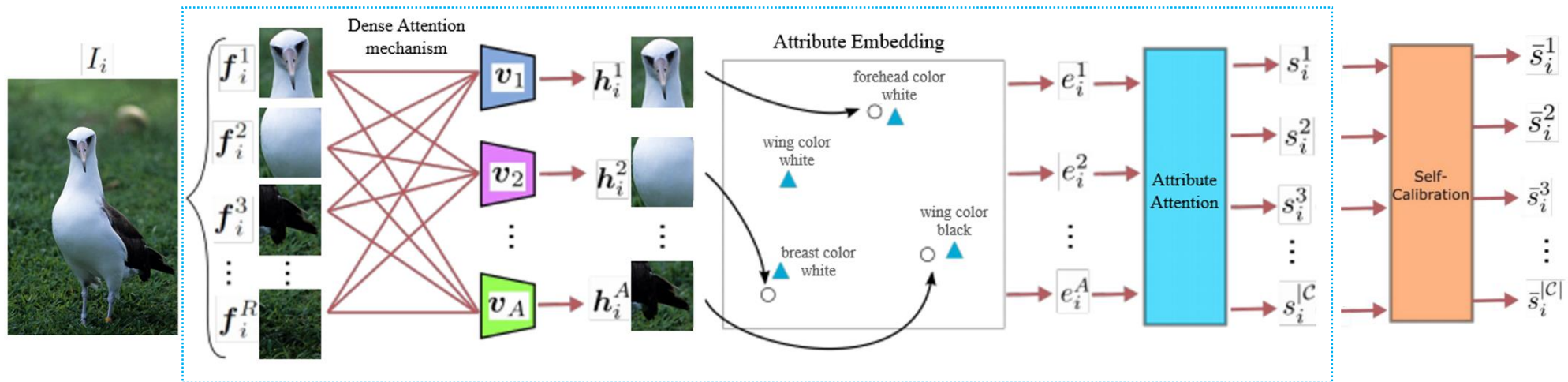
Very Large seabird with narrow, pointed wings and distinctive looping flight style; rarely flaps wings. Underwings are patchy white with variable amounts of dark. Grayish smudge on face.



Juvenile

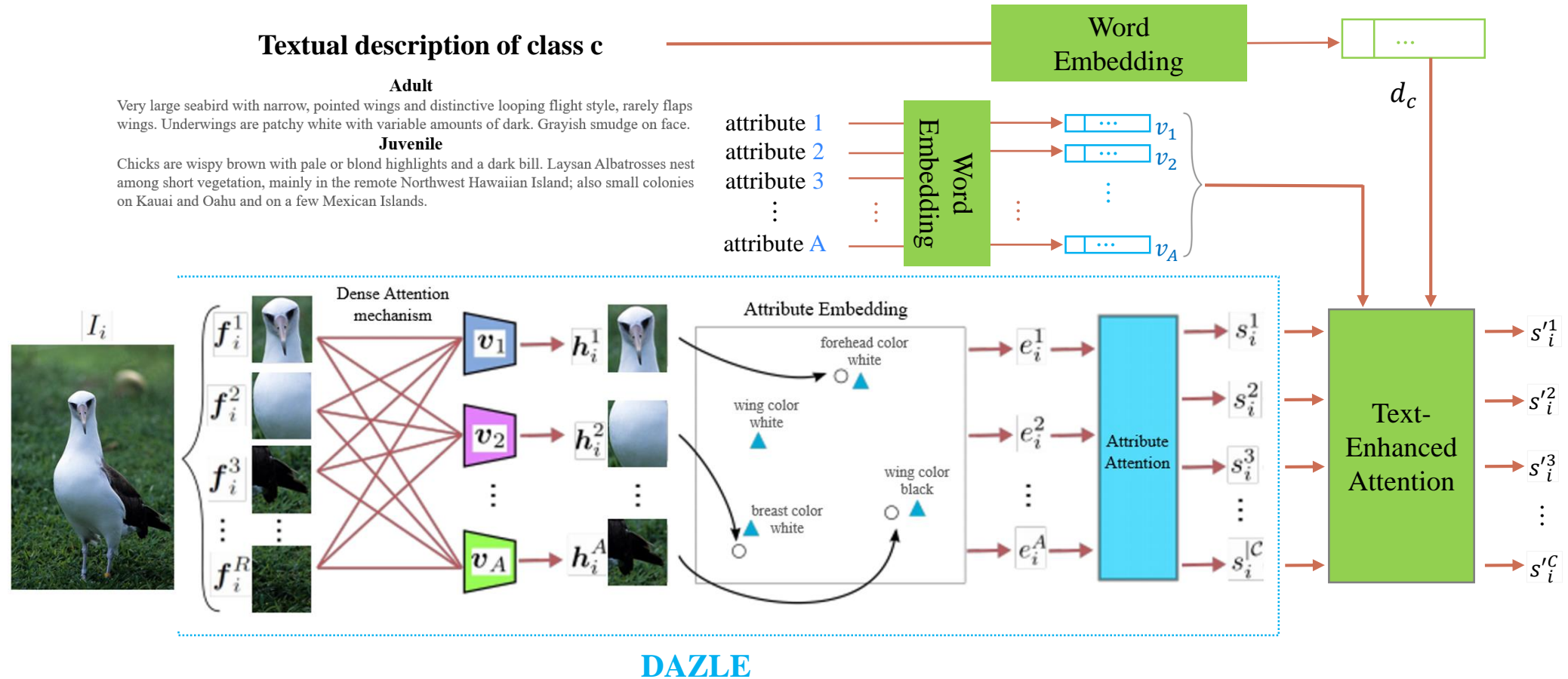
Chicks are wispy brown with pale or blond highlights and a dark bill. Laysan Albatrosses nest among short vegetation, mainly in the remote Northwest Hawaiian Islands; also small colonies on Kauai and Oahu and on a few Mexican Islands.

# Text-enhanced attribute-based attention



DAZLE

# Text-enhanced attribute-based attention



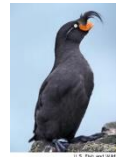


# Datasets

Dataset	attributes	seen (val) / unseen classes	training / testing samples
CUB	312	100 (50) / 50	7,057 / 4,731
AWA2	85	27 (13) / 10	23,527 / 13,795



AWA2  
(Coarse-grained)



CUB  
(Fine-grained)

# Results

Method	CUB			AWA2		
	$acc_s$	$acc_u$	H	$acc_s$	$acc_u$	H
DAZLE (dense attention)	57.6	42.4	48.8	72.1	46.8	56.8
Ours (text-enhanced attention)	<b>60.9</b>	<b>46.5</b>	<b>52.6</b>	<b>72.8</b>	<b>59.9</b>	<b>65.8</b>

- Our approach improves DAZLE in both datasets.
  - On CUB, the accuracy rates of recognizing seen and unseen classes are boosted by 2.4% and 4.1%, respectively. The harmonic mean is increased by 3.8%.
  - A performance gap for the unseen classes on AWA2 is also observed. The harmonic mean is increased by 9%.



# Conclusions

- We improve DAZLE by devising a *text-enhanced* attribute-based attention mechanism that guides the feature extraction from most relevant image regions for important attributes.
- By leveraging *textual descriptions*, the proposed method adapts well to unseen classes at inference.

More information:

<http://www.csie.ntnu.edu.tw/~myeh>