



# Zero-shot learning

Soichiro Tanaka, Inspection室

2019/12/26



## Slide Title

---

- What is zero-shot learning?
- Why do we need zero-shot learning?
- How does it work?
- Experiment
- Conclusion

# What is zero-shot learning?



He has seen a horse before.



He has never seen a zebra.

# What is zero-shot learning?

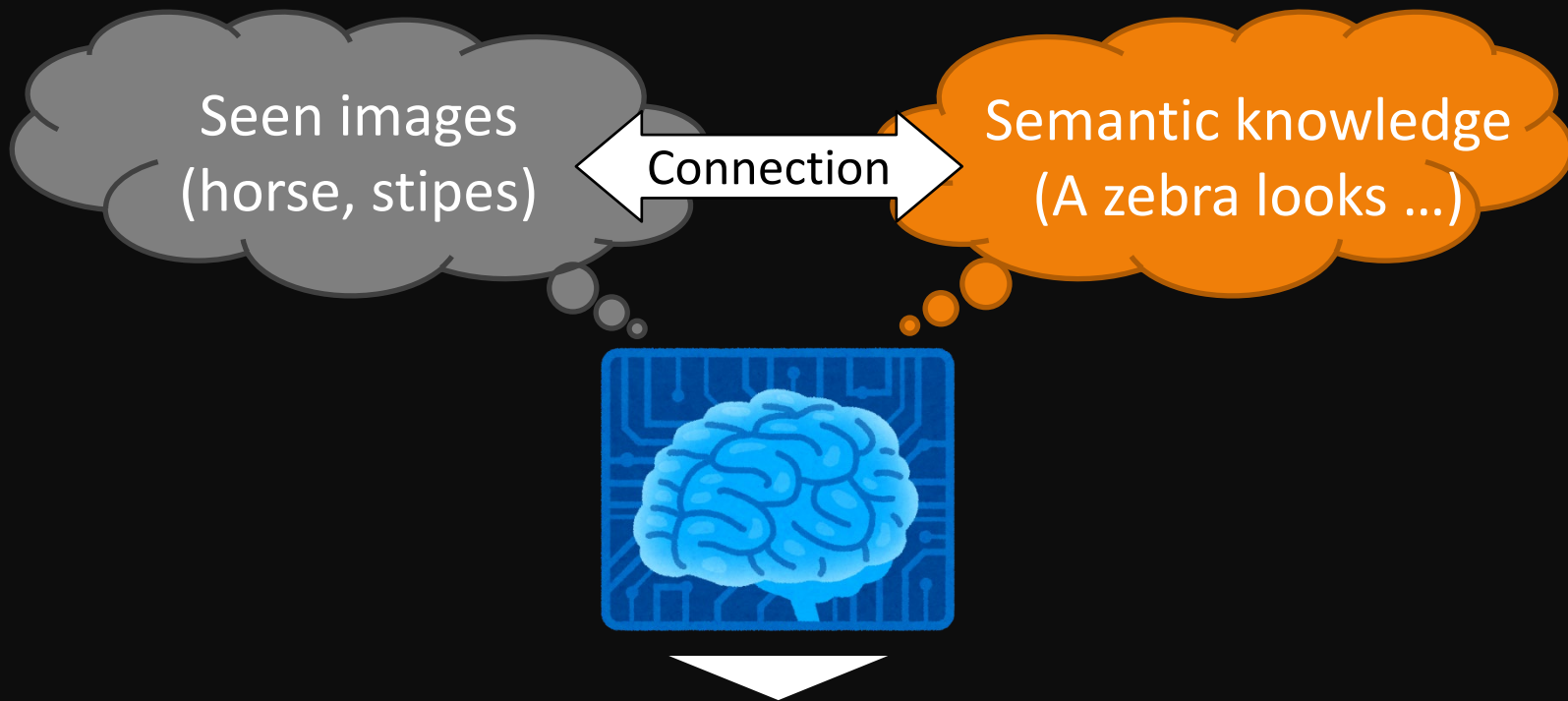


But he has read somewhere  
“A zebra looks similar to a horse,  
but has black-and-white stripes.”



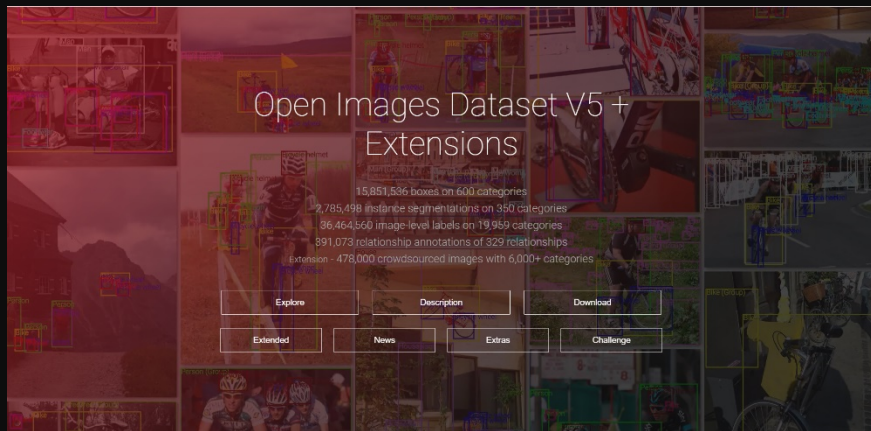
# What is zero-shot learning?

## Motivation



# Why do we need zero-shot learning?

Because of vast and growing number of categories.



15,851,536 boxes on 600 categories  
2,785,498 instance segmentations on 350 categories  
36,464,560 image-level labels on **19,959 categories**  
391,073 relationship annotations of 329 relationships

Adding new category to DL model requires a lot of labor.

- ✓ Collecting **enough** data for each class
- ✓ **Annotation** cost is high
- ✓ Change the output layer and **train again**
- ✓ More **time** or **computational power** with more data

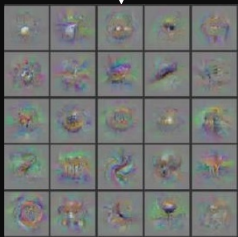
Zero-shot learning can solve these problems maybe.

# How does it work?

## Learning phase



feature extraction



word vector

Word embedding space

Learn how to map



“horse”

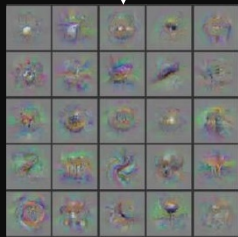
word2vec

word vector

- eat grasses
- run fast
- hunted
- brown
- ...

# How does it work?

## Prediction phase



- word vector
- horse-like
  - black-and-white
  - stripes
  - ...

Word embedding space

Search nearest vector

Generate word from vector



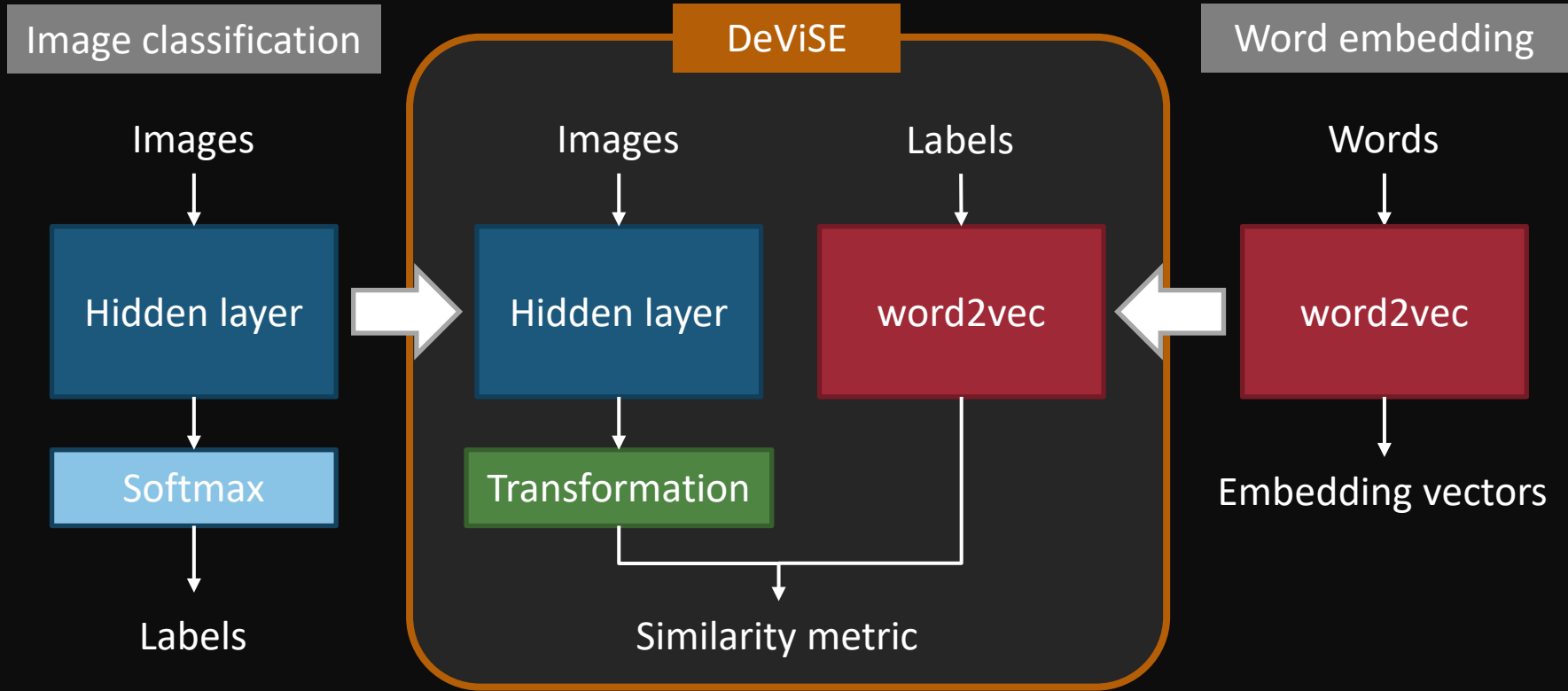
“zebra”

word2vec



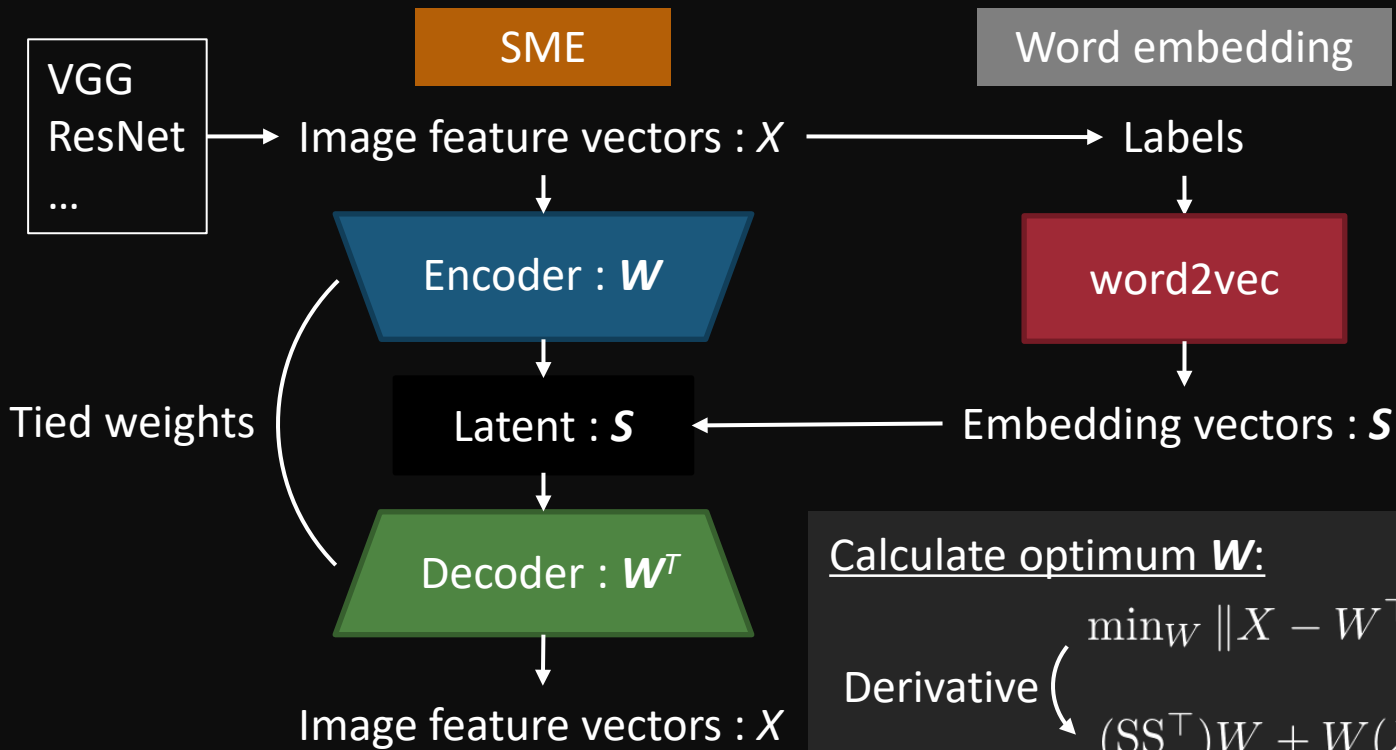
# How does it work?

## Deep Visual-Semantic Embedding(DeViSE) NIPS2013



# How does it work?

## Semantic AutoEncoder(SAE) CVPR2017



Calculate optimum  $W$ :

$$\min_W \|X - W^T S\|_F^2 + \lambda \|WX - S\|_F^2$$

Derivative

$$\underbrace{(SS^T)}_{const} W + W \underbrace{(\lambda X X^T)}_{const} = \underbrace{(1 + \lambda) S X^T}_{const}$$

# Experiment

## Animal classification

Dataset : A Large-scale Attribute Dataset for Zero-shot Learning



### Bear

Color: is white: True  
Limb: has short legs: True  
Behaviour: can swim: True  
Habit: lives in groups: False



### Lychee

Size: is big: False  
Shape: is globular: True  
Edibility: has nutlets: True  
Medicinal property: is mild: True



### Fighter

Parts: has a jet engine: True  
Color: is green : False  
Safety: is dangerous : True  
Power consumes: wind power : False



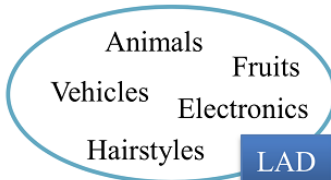
### Mobile Phone

Parts: has a battery: True  
Shape: is flat: True  
Function: can photograph: True  
Aim: is for cleaning: False



### Bob Hair

Color: is brown: True  
Color: is black: False  
Fitness: fits people with earring: False  
Feeling: is cute: True



ID	class	ID	class 11
0	ant	25	jellyfish
1	bear	26	koala
2	butterfly	27	leopard
3	carp	28	lion
4	cat	29	monkey
5	catfish	30	ostrich
6	cattle	31	panda
7	chicken	32	parrot
8	coral	33	platypus
9	deer	34	rabbit
10	dog	35	rat
11	dolphin	36	rhinoceros
12	donkey	37	shark
13	dragonfly	38	snail
14	duck	39	sparrow
15	eagle	40	starfish
16	elephant	41	sturgeon
17	fox	42	swallow
18	giraffe	43	swan
19	goat	44	tiger
20	goose	45	turtle
21	gorilla	46	whale
22	hedgehog	47	wolf
23	honeybee	48	woodpecker
24	horse		

Training: 9048 images

Validation: 3878 images

Unseen classes : zebra, chimpanzee

# Experiment

## Network architecture

VGG16 with two branches

- Mapping from image features to word embedding space (300 dim)
- Classification (49 categories)

## Settings

Optimizer: adam

Learning rate:  $1e-3 \rightarrow 1e-6$

Epochs: 1000

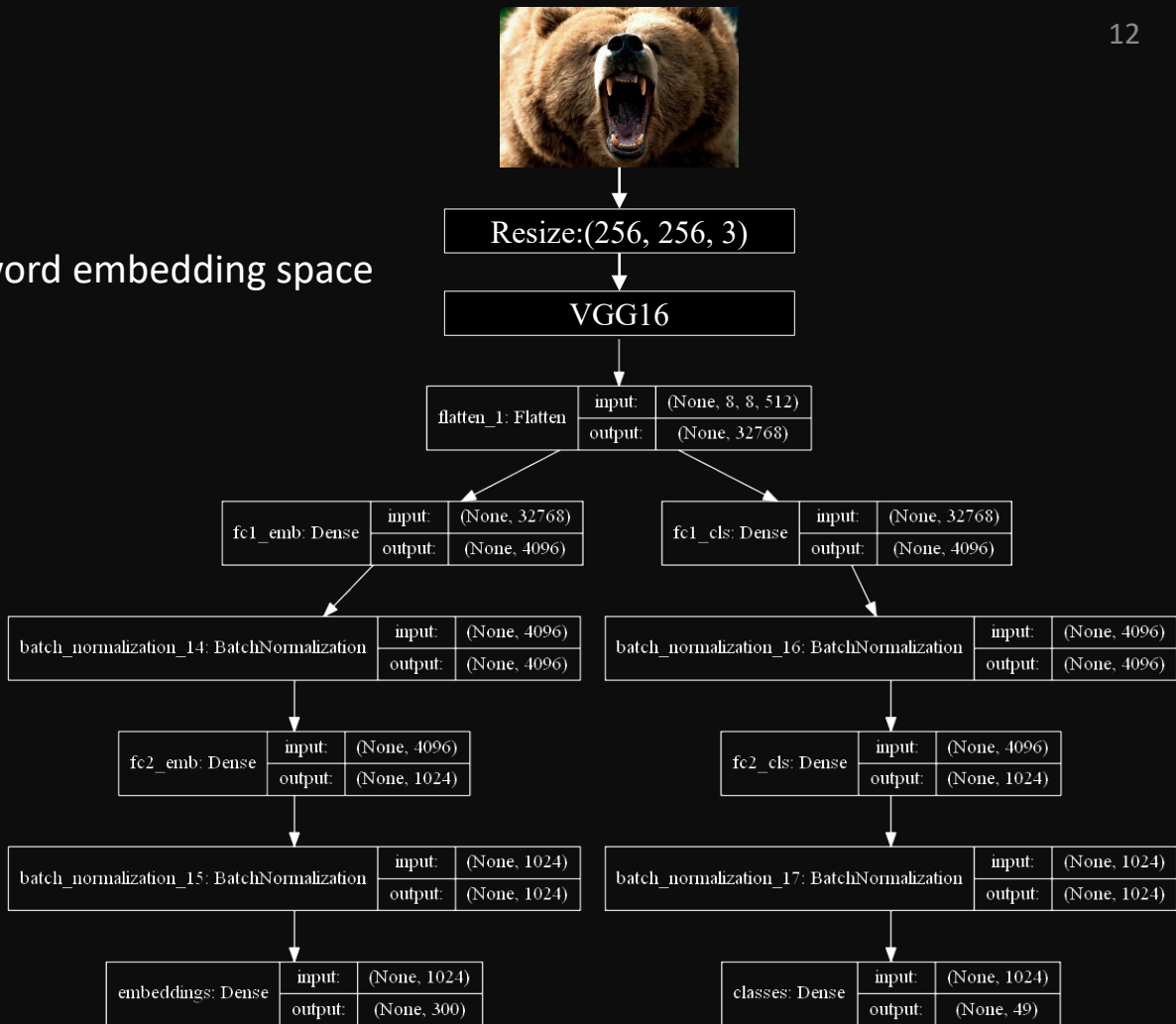
Batch size: 32

Input shape: (256, 256, 3)

Loss1: cosine proximity










Loss2: categorical crossentropy

Loss weights: 'loss1': 0.7, 'loss2': 0.3



# Experiment

## Prediction of unseen classes

Input	Output	Reliability	Input	Output	Reliability	Input	Output	Reliability
	woodpecker	0.857601		turtle	0.886721		goat	0.818994
	<b>squirrel</b>	<b>0.707087</b>		<b>alligator</b>	<b>0.754621</b>		<b>cow</b>	<b>0.737375</b>
	spotbill	0.704446		tortoise	0.736731		dog	0.719159
	liventer	0.698276		turtles	0.719899		pig	0.716362
	gymnogene	0.696299		<b>crocodile</b>	<b>0.716109</b>		rabbit	0.707818
	rhinoceros	0.985257		<b>chicken</b>	<b>0.847994</b>		butterfly	0.778441
	rhinoceroses	0.719658		goat	0.815152		parrot	0.687232
	dicerorhinus	0.695453		cow	0.764047		turtle	0.647621
	<b>hippopotamus</b>	<b>0.692935</b>		pig	0.760932		<b>frog</b>	<b>0.636198</b>
	chousingha	0.683939		chickens	0.7592		spotbill	0.632529
	tiger	0.896356		starfish	0.945515		gorilla	0.993781
	leopard	0.765723		coral	0.776223		ape	0.655898
	<b>hyena</b>	<b>0.69587</b>		<b>urchins</b>	<b>0.732707</b>		<b>chimpanzee</b>	<b>0.655245</b>
	rhinoceros	0.686366		<b>urchin</b>	<b>0.7239</b>		monkey	0.645724
	elephant	0.67506		corals	0.722444		<b>chimp</b>	<b>0.639019</b>

# Experiment

Zebra's stripes are overwhelming!



leopard	0.954025
profelis	0.627406
panthera	0.612081
chousingha	0.597279
chowsingha	0.590462



leopard	0.902383
giraffe	0.72913
profelis	0.713392
hyena	0.705404
rhinoceros	0.696642



giraffe	0.794636
leopard	0.784001
animazia	0.716992
hyena	0.715736
profelis	0.706123



leopard	0.945774
tiger	0.781351
panthera	0.675595
hyena	0.647742
giraffe	0.64211

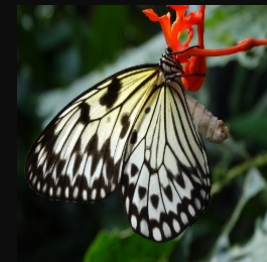
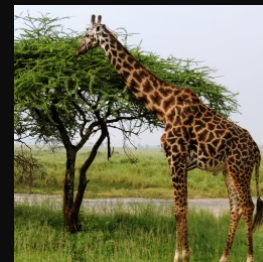
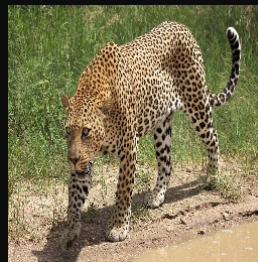


butterfly	0.773512
leopard	0.771586
tiger	0.649144
animazia	0.593183
profelis	0.576424



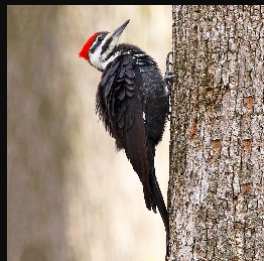
butterfly	0.868839
tiger	0.685825
leopard	0.618752
medley	0.594179
swallowtail	0.572934

Zebra images are classified into around leopard, giraffe, butterfly because of **weak distributed representation** power **except for stripes**.



# Experiment

## Woodpecker vs Wood



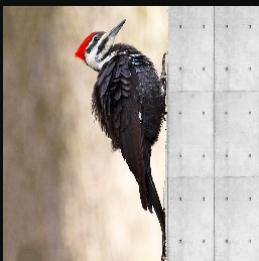
<b>woodpecker</b>	<b>0.992235</b>
dendrocopos	0.71069
dryocopus	0.691981
leucotos	0.684678
campethera	0.683521



<b>woodpecker</b>	<b>0.954602</b>
bellied	0.693659
parrot	0.691737
dendrocopos	0.687798
pileated	0.678356



<b>woodpecker</b>	<b>0.991888</b>
dendrocopos	0.718132
dryocopus	0.696118
campethera	0.685811
leucotos	0.682005



squirrel	0.745027
profelis	0.722475
animazia	0.718295
parrot	0.714114
turtle	0.710388

Misrecognize tree as woodpecker because woodpecker is always with tree.

→ ZL maps predominantly the feature of tree into word embedding space.

### Solutions

- ✓ Add images of other animals with tree to train data.
- ✓ Attract attention to animals (remove background, attention network, etc.)

# Conclusion

## Result

Succeeded in predicting unseen categories very rarely.

## Problems

Feature extraction fails to cover all important features because...

- texture is overwhelming
- background attracts much attention

## Feature work

- ✓ Study attention networks.
- ✓ Knowledge Graph greatly enhances the performance.
  - Zero-shot Recognition via Semantic Embeddings and Knowledge Graphs (CVPR 2018)
- ✓ Try object detection and semantic segmentation with ZL.
- ✓ Apply to actual projects in the distant future.



# Unsuitable problem

