



Beyond Seen Primitive Concepts and Attribute-Object Compositional Learning

Nirat Saini Khoi Pham Abhinav Shrivastava

University of Maryland, College Park

Compositional Zero-shot Learning (CZSL)

Input images

Output labels

Goal: To compose seen attributes and objects for recognizing unseen compositions.

Method: Semantic correlations in linguistic embedding space help recognizing unseen compositions.

Assumption: All objects and attributes are individually seen during training.

Practical Settings: All objects and attributes are not present in training set. How to learn composition for unseen attributes and objects?

Proposed method is built on top of OADis¹ (baseline for CZSL)

- Images and labels are embedded together (Vision+Language : Pair Embedding Space).
- Visual features are separated for attributes and objects (Object-Attribute Disentanglement).
- Composition classification is a separate head (Label Embedder).
- Each embedding space (attribute, object and pair) is regularized using linguistic losses.
- Neighborhood Expansion Loss transfers knowledge from seen to unseen concepts and compositions.

Pipeline

Input image

Pair Embedding Space

Composition Labels

External Knowledge

Neighborhood Expansion Loss (NEL) = $H(y, C) = \sum_{m=1}^M -y_m \log(C_m)$

Neighborhood Expansion Loss is combination of:

- Cross Entropy Loss (CE): $H(y, C) = \sum_{m=1}^M -y_m \log(C_m)$
- Label Smoothing (LS): Avoids negative bias towards unseen labels
- Label Propagation (LP): Transfer from seen to unseen labels (from External Knowledge Source)

Open Vocabulary Compositional Zero-shot Learning (OV-CZSL)

Semantic Correlations in Language

External Knowledge Sources

Goal: To recognize attributes, objects and their compositions beyond seen vocabulary.

Training: set of compositions of seen attributes and objects

Testing: 4 splits of compositions of seen and unseen attributes and objects

Method: Use External Knowledge source and transfer knowledge from seen to unseen attributes, objects and their compositions.

In-the-wild generalization for Open Vocabulary Compositions. Open Vocabulary Compositional Zero-shot Learning: ✓

Test splits for CZSL and OV-CZSL

	Attributes	Objects
Seen	Seen	Unseen
Unseen	Unseen	Unseen
peeled lemon	sliced grapefruit	?
peeled ~ sliced	lemon ~ grapefruit	
peeled lemon ~ sliced grapefruit		

Compositional Zero-shot Learning Task: ★

Quantitative Results on Benchmark splits

Zero-shot Learning Baselines for Seen and Unseen Attribute and Object Accuracy

NEL as plug-and-play module improves Unseen Composition Accuracy for baselines

- Experiments are done on a new benchmark split for OV-CZSL for existing datasets.
- We report Top@1 AUC for MIT-states and CGQA, and Top@3 and 5 AUC for VAW-CZSL.
- AUC (%) is computed between seen and unseen compositions with different bias terms and chosen where HM is maximum.

Results

Attribute Cluster Performance

Object Cluster Performance

Top 3 predictions for Images