



Generating Visual Explanations

Lisa Anne Hendricks¹, Zeynep Akata², Marcus Rohrbach¹,
Jeff Donahue¹, Bernt Schiele², Trevor Darrell¹

¹University of California Berkeley, ²Max Planck Institute for Informatics



Motivation

- Explanations are important for understanding and interacting with intelligent systems

Visual Explanations

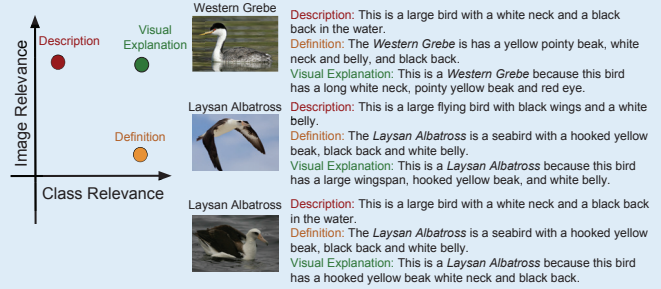
- Visual explanations must be image relevant and class consistent

Justification vs. Introspection

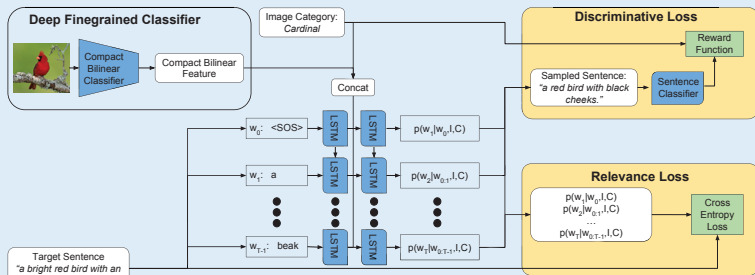
- Justification* explanations detail why a prediction is compatible with visual evidence
- Introspective* explanations detail the internal mechanisms of a model
- Here, concentrate on justifications as they are useful for non-experts

Method:

- Introduce novel discriminative loss based on REINFORCE [1]



Model



Relevance Loss

Encourages sentences to be **image relevant** by minimizing cross entropy between ground truth and predicted words.

$$L_R = \sum_{t=0}^{T-1} \log p(w_{t+1} | w_{0:t}, I, C)$$

where I is an image feature, w_t is a ground truth word, and C is a class label.

Discriminative Loss

Encourages sentences to be **class discriminative** by maximizing a reward, R_D , which measures class discriminativeness of a sampled sentence \tilde{w} .

$$L_D = -E_{\tilde{w} \sim p(w)} [R_D(\tilde{w})]$$

Estimate L_D with Monte Carlo sampling because computing an expectation over descriptions is intractable.

Use REINFORCE [1] to estimate the expected gradient of L_D :

$$\nabla E_{\tilde{w} \sim p(\tilde{w})} [R_D(\tilde{w})] = E_{\tilde{w} \sim p(\tilde{w})} [R_D(\tilde{w}) \nabla_W \log p(\tilde{w})]$$

Data

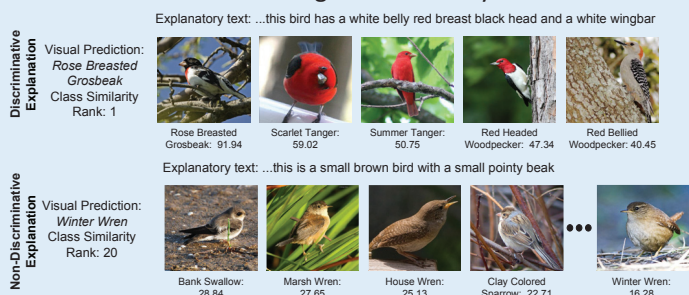
- Use finegrained descriptions collected in [2]
- 200 bird classes, > 11k images, 5 captions/image
- Finegrained descriptions more class informative than attributes.



Ground truth descriptions:
 ...there is a **black bird** with a **black beak** and **red eyes**.
 ...this bird has a **large, black, curved bill**, **black tarsuses** and **feet**, and a **black throat, breast**, and **belly**.
 ...**angular beak** and **red eye**, this otherwise nondescript black bird has **alarmingly large feet**.

Metrics

- Measure **image relevance** using CIDEr.
- Measure **class relevance** using class similarity rank.



Results

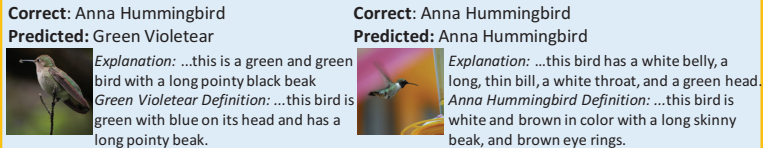
Explanations are more **class relevant** than descriptions.



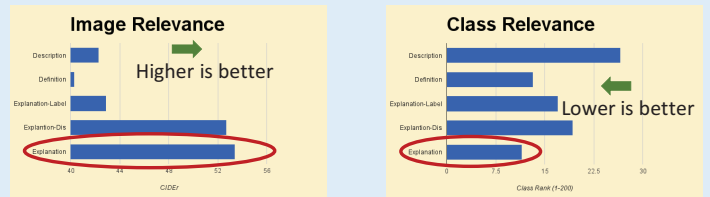
Explanations are more **image relevant** than definitions.



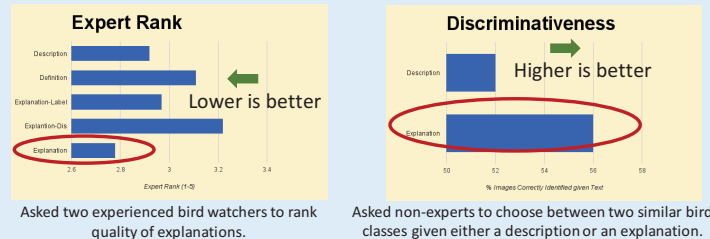
Explanations for incorrect classification decisions.



Quantitative Evaluations



Human Evaluations



[1] Williams, Ronald J. "Simple statistical gradient-following algorithms for connectionist reinforcement learning." *Machine learning* 8.3-4 1992.

[2] Reed, Scott, et al. "Learning Deep Representations of Fine-Grained Visual Descriptions." *CVPR 2016*.